



# ScRNASeqSnap: Single cell RNA Seq Snap workflow

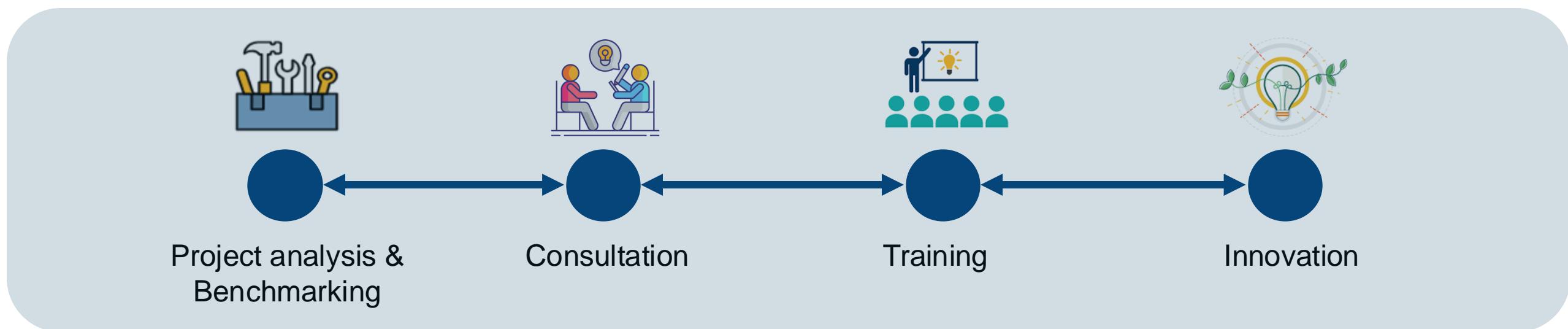
Antonia (Tonia) Chroni, PhD  
Senior Bioinformatics Research Scientist  
DNB Bioinformatics Core



# Bioinformatics core, Department of Developmental Neurobiology



**Providing advanced bioinformatic services for investigators to leverage omics data**



**Cody Alexander Ramirez, PhD**  
Senior Bioinformatics Research Scientist  
Core Director  
Boston, Massachusetts



**Antonia Chroni, PhD**  
Senior Bioinformatics Research Scientist  
New York, New York



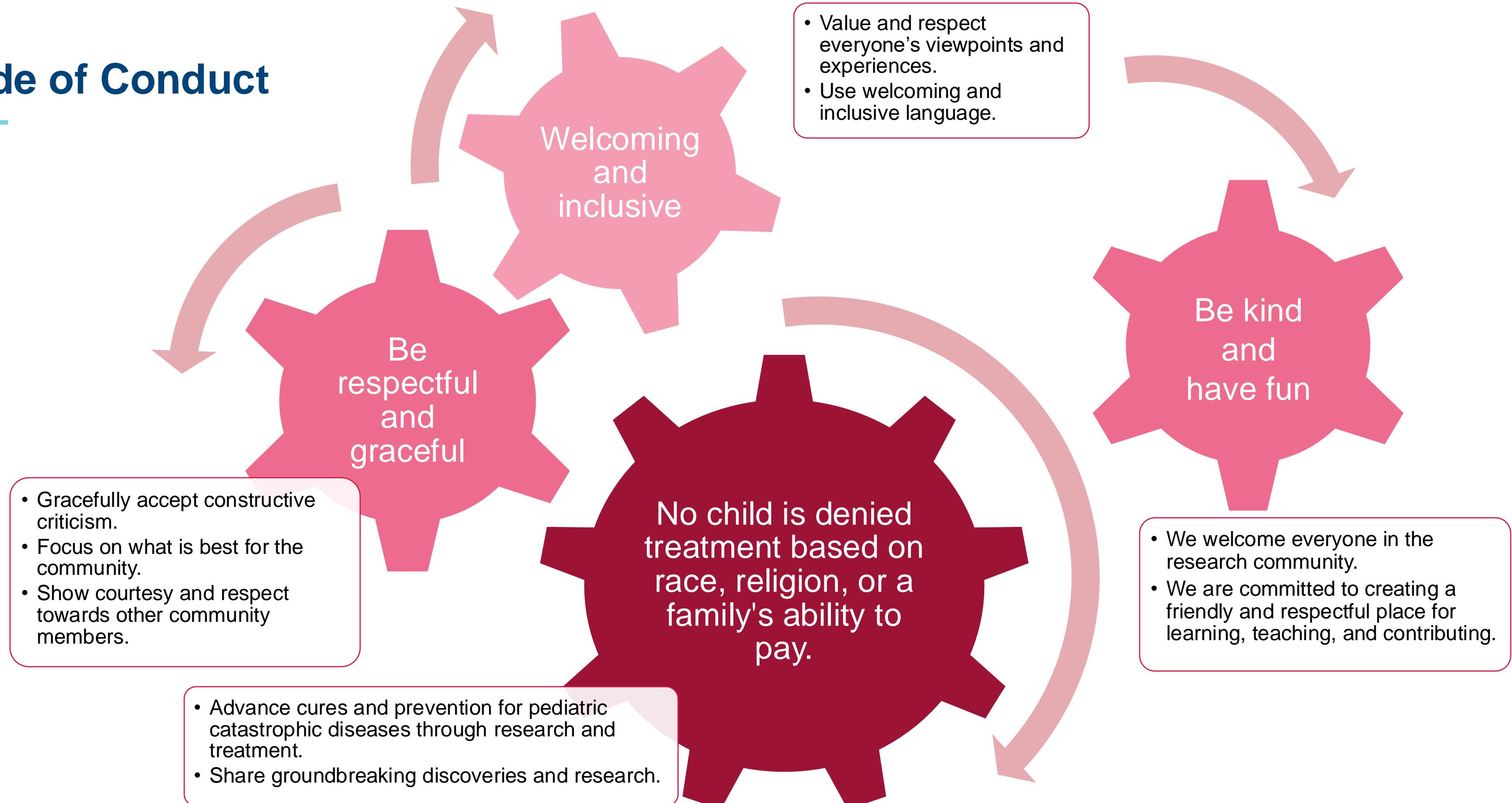
**Sharon Freshour, PhD**  
Bioinformatics Research Scientist  
St. Louis, Missouri



**Asha Jacob Jannu, PhD**  
Bioinformatics Research Scientist  
Indianapolis, Indiana



# Code of Conduct



# Learning Objectives

1. Overview of Tools and Methods in the `ScRNASeqSnap` workflow
2. Reviewing Results, Plots, and Reports



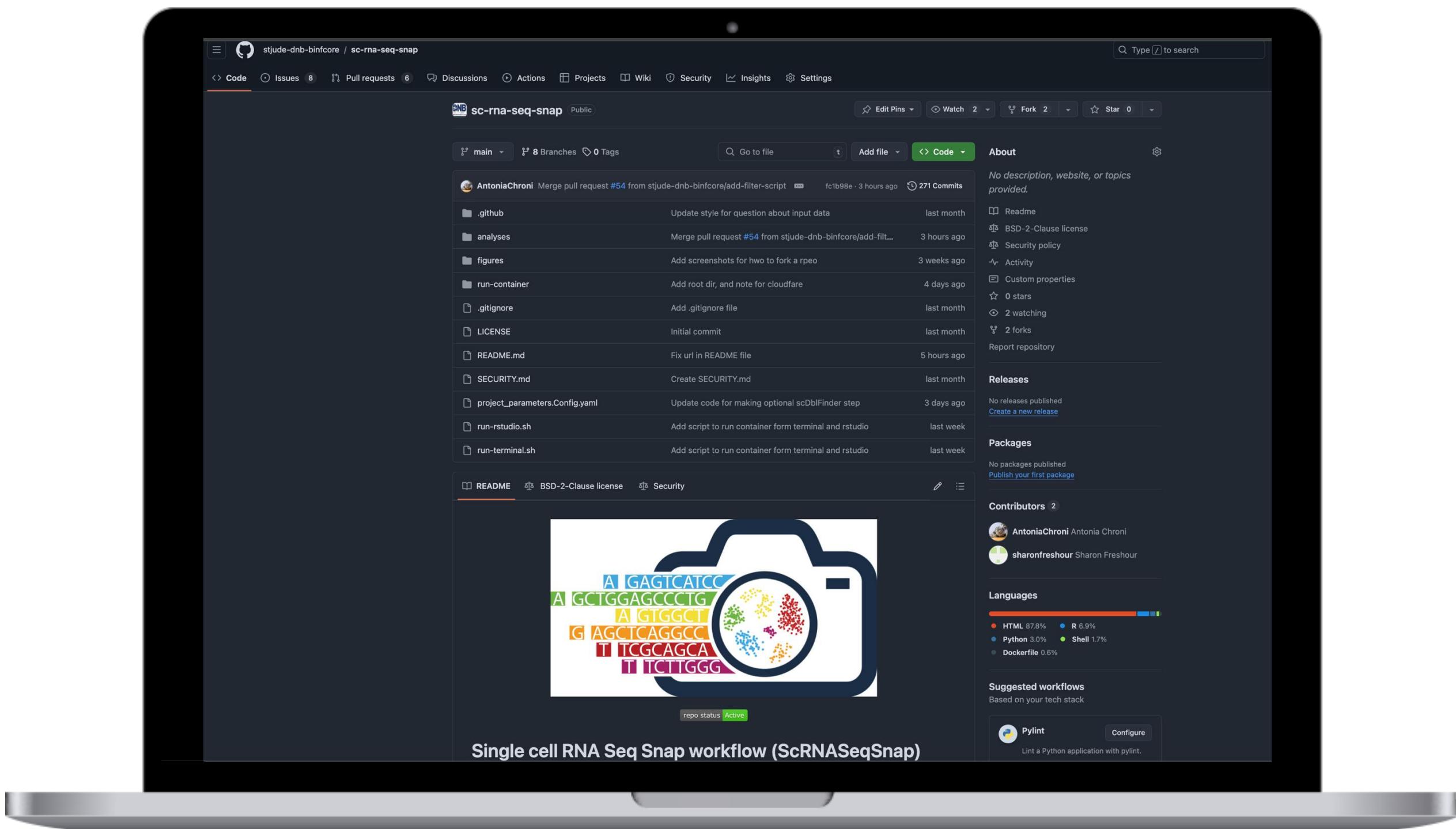
# Single cell RNA Seq Snap workflow (ScRNASeqSnap)

---

- 🚀 Automated suite of analysis modules for upstream and downstream analyses
- 🧬 sc/snRNA data analysis from 10X sequencing technology for non-hashed experiments
- 🧣 Human, mouse, and dual index genome experiments
- 📝 Thorough and step-by-step documentation on how to use each analysis module
- 👉 Reproducible and replicable results via Docker/singularity
- 🌐 Publicly available at [stjude-dnb-binfcore/sc-rna-seq-snap](https://stjude-dnb-binfcore/sc-rna-seq-snap)



# Single cell RNA Seq Snap workflow (ScRNASeqSnap)



<https://github.com/stjude-dnb-bincore/sc-rna-seq-snap>





# How to access the code

The screenshot shows a GitHub repository page for `stjude-dnb-bincore/sc-rna-seq-snap`. The main heading is "To access the code in this repository:". Below it, there's a "Clone option" section with the instruction "1. Clone the repository" and a command-line example:

```
git clone https://github.com/stjude-dnb-bincore/sc-rna-seq-snap.git
```

Below this, there's a "Fork option" section with the instruction "1. Fork the repository on your own account from the main page of the `stjude-dnb-bincore/sc-rna-seq-snap` by clicking the "Fork" button". It shows a screenshot of the repository's main page with the "Fork" button highlighted.

Further down, there's another step: "2. Change the name if you like, but probably not; click "Create fork"" followed by a screenshot of the "Create a new fork" dialog box.





# How to access the code

The screenshot shows a GitHub repository page for `stjude-dnb-bincore/sc-rna-seq-snap`. The main heading is "To access the code in this repository:". A red box highlights the "Clone option" section, which contains the instruction "1. Clone the repository" and the command:

```
git clone https://github.com/stjude-dnb-bincore/sc-rna-seq-snap.git
```

Below this is the "Fork option" section, which contains the instruction "1. Fork the repository on your own account from the main page of the `stjude-dnb-bincore/sc-rna-seq-snap` by clicking the "Fork" button". It shows a screenshot of the repository's main page with the "Fork" button highlighted. The next step, "2. Change the name if you like, but probably not; click "Create fork"" is shown below, with another screenshot of the "Create a new fork" dialog.





# How to access the code

To access the code in this repository:

**Clone option**

1. Clone the repository

```
git clone https://github.com/stjude-dnb-bincore/sc-rna-seq-snap.git
```

**Fork option**

1. Fork the repository on your own account from the main page of the `stjude-dnb-bincore/sc-rna-seq-snap` by clicking the "Fork" button

**2. Change the name if you like, but probably not; click "Create fork"**



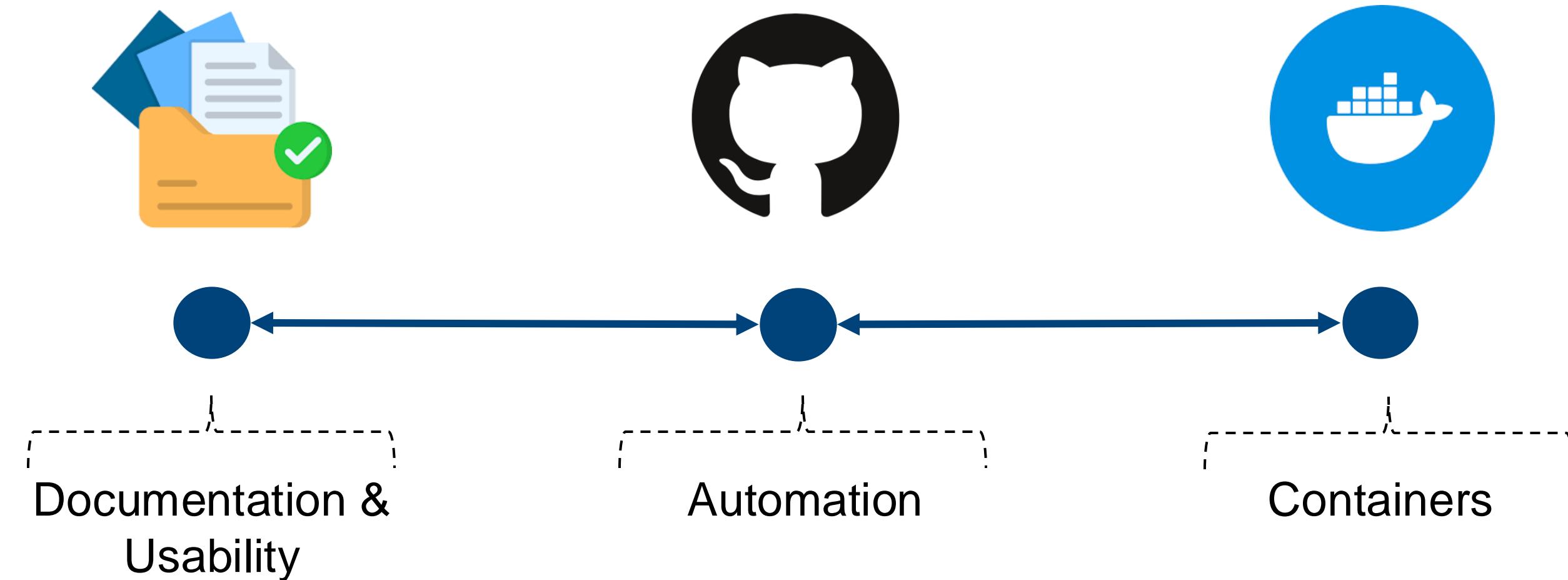


# ScRNASeqSnap is automated, reproducible and replicable

---



# ScRNASeqSnap is automated, reproducible and replicable

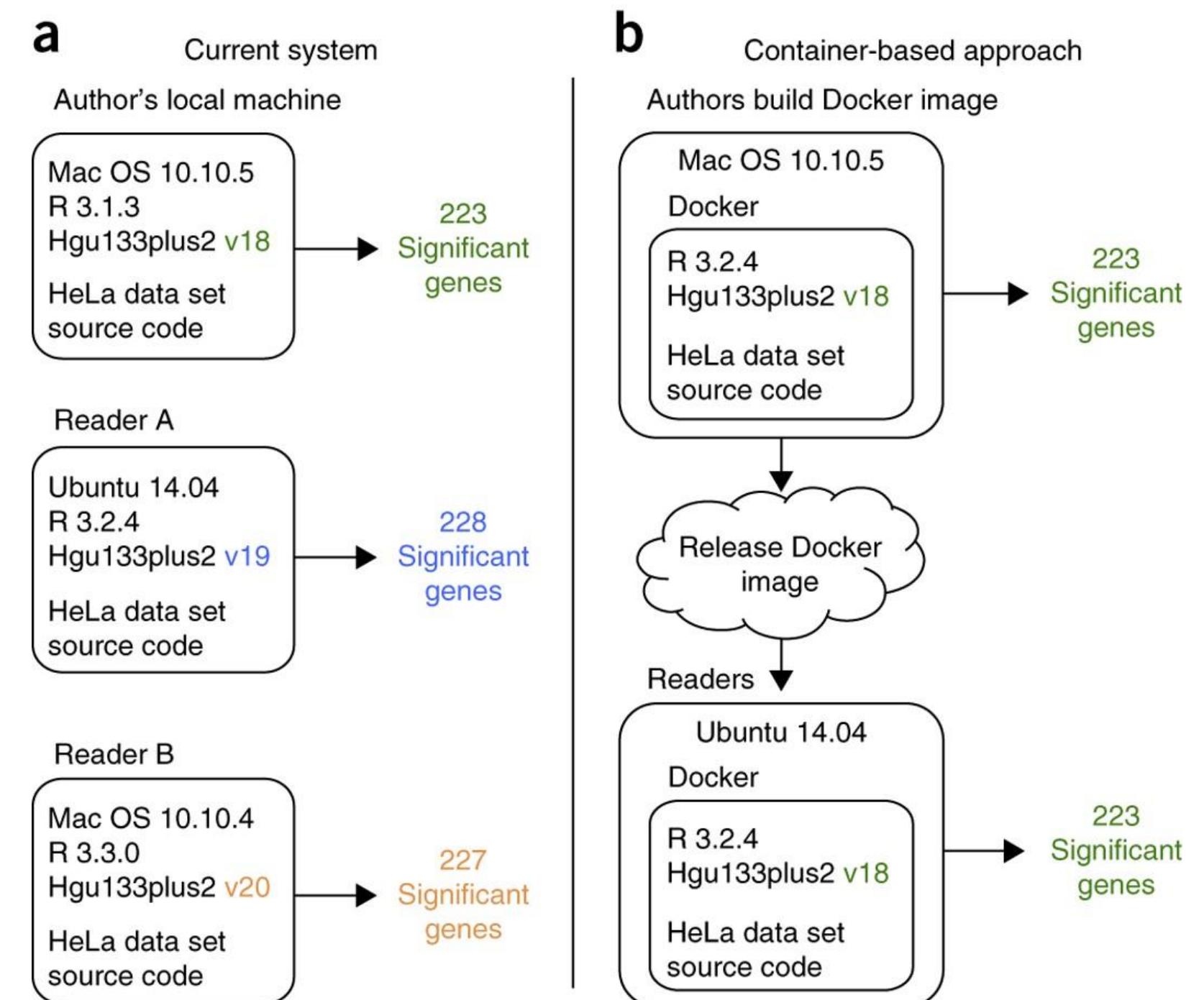


Slide from the "["Automation and Reproducibility in Computational Biology"](#) course



# ScRNASeqSnap is automated, reproducible and replicable

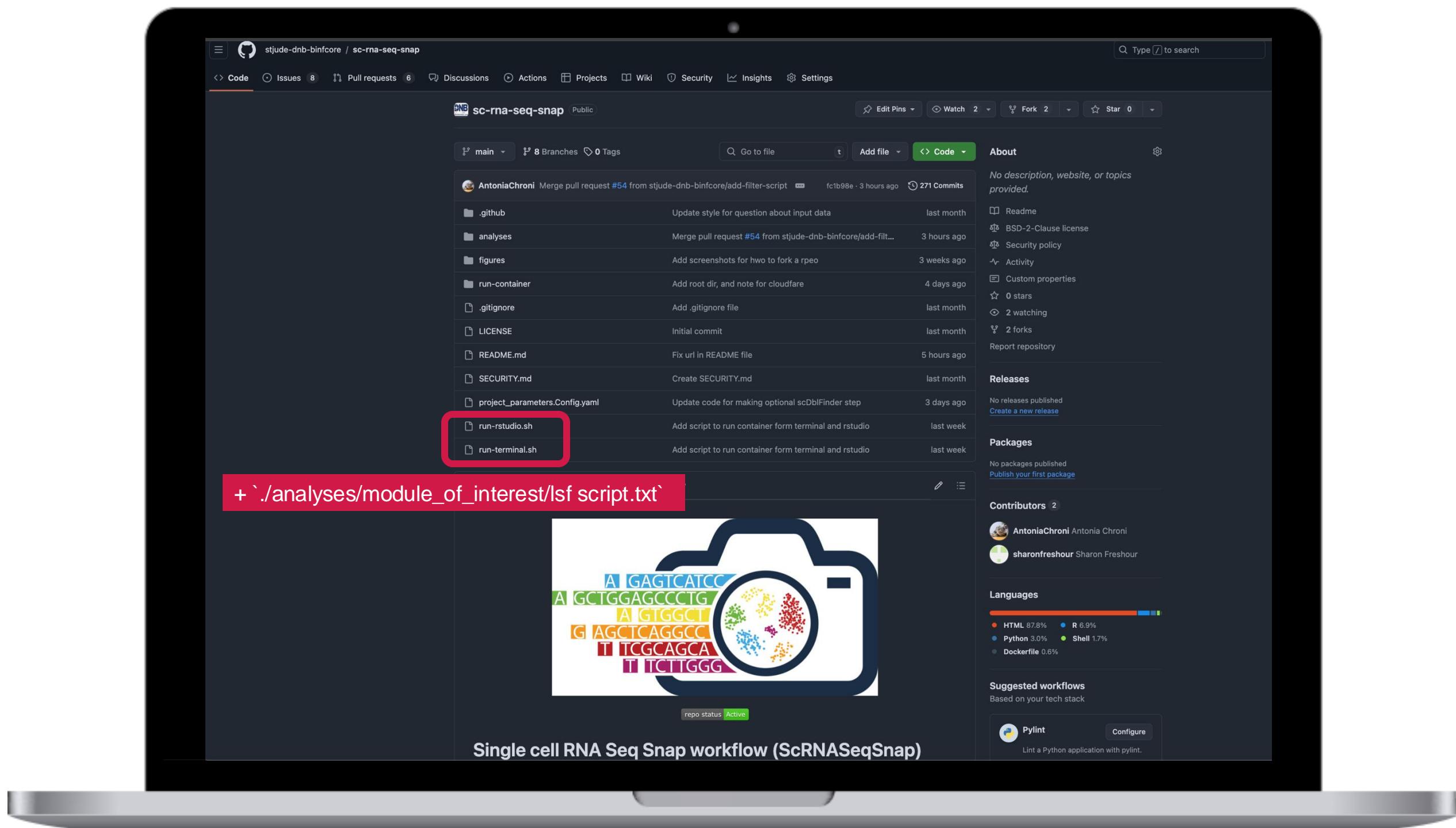
- **Dockerfile/rstudio\_r\_4.4.0\_seurat\_4.4.0.def**
- Docker and Singularity
  - “Containers” that include everything from the operating system up
  - Run one OS inside another, with all the things frozen to particular versions.



Slide from the "[Automation and Reproducibility in Computational Biology](#)" course



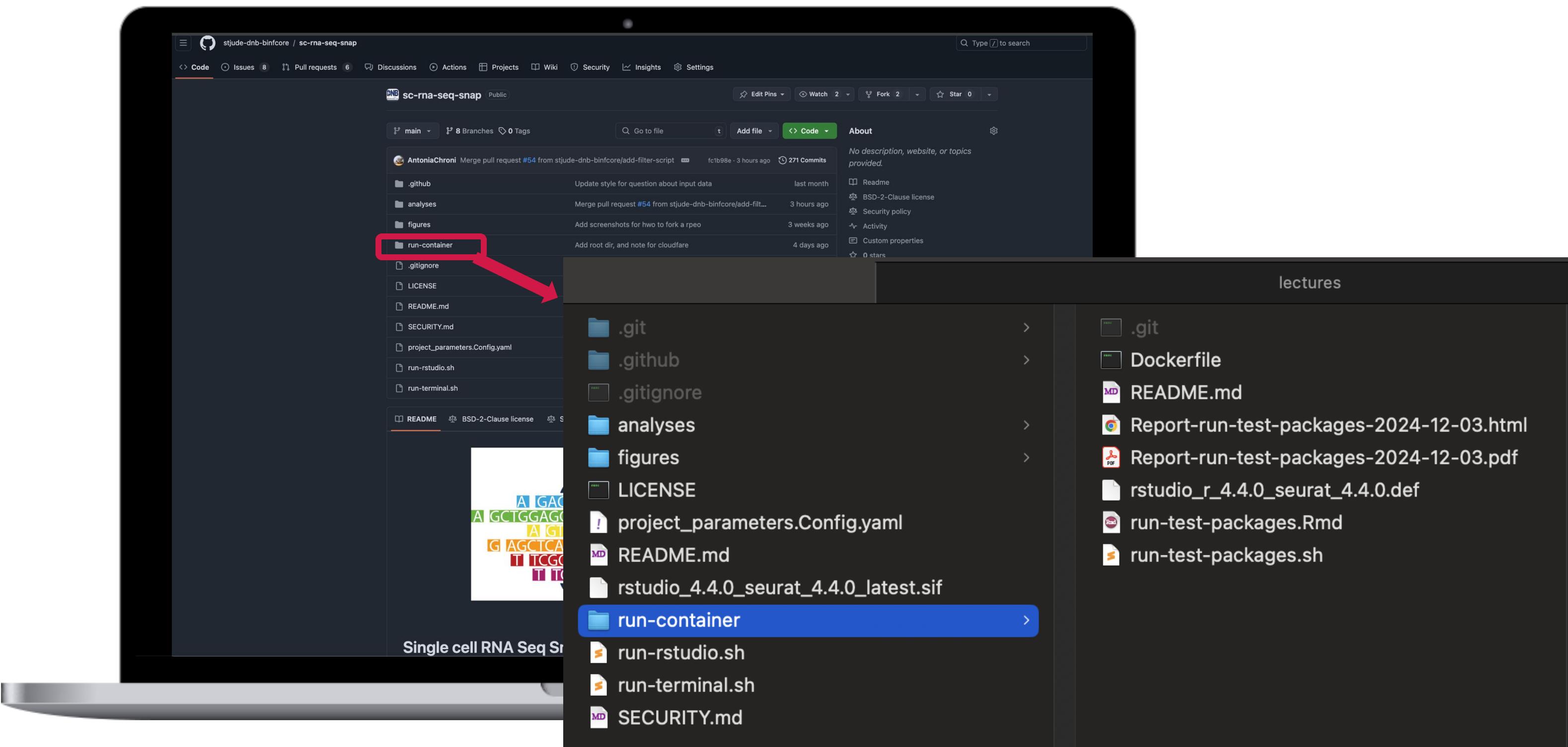
# ScRNASeqSnap is automated, reproducible and replicable



<https://github.com/stjude-dnb-bincore/sc-rna-seq-snap>



# ScRNASeqSnap is automated, reproducible and replicable





# Repository organization and project tidiness

---

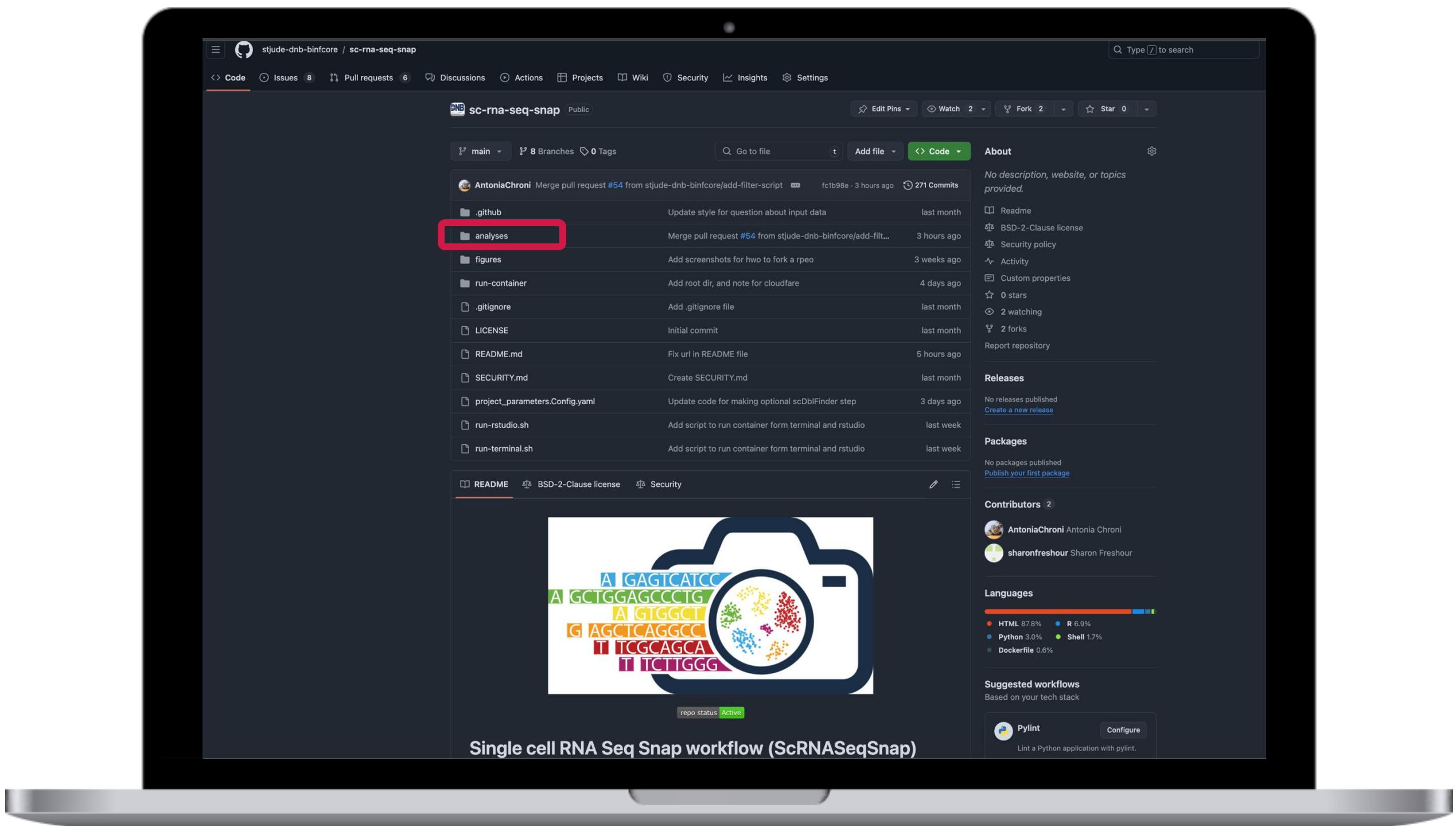
- **Use Folders/Directories**
- **Keep separate projects separated**
- **Separate sections for units within a project**
  - analyses/<module\_name>
    - code
    - plots
    - results
  - data
  - figures
  - LICENSE
  - project\_parameters.Config.yaml
  - SECURITY.md
  - ...
- **Documentation throughout:** Describe what files do and how they are organized/folder.

Slide from the "[Automation and Reproducibility in Computational Biology](#)" course





# Repository organization and project tidiness



<https://github.com/stjude-dnb-bincore/sc-rna-seq-snap>





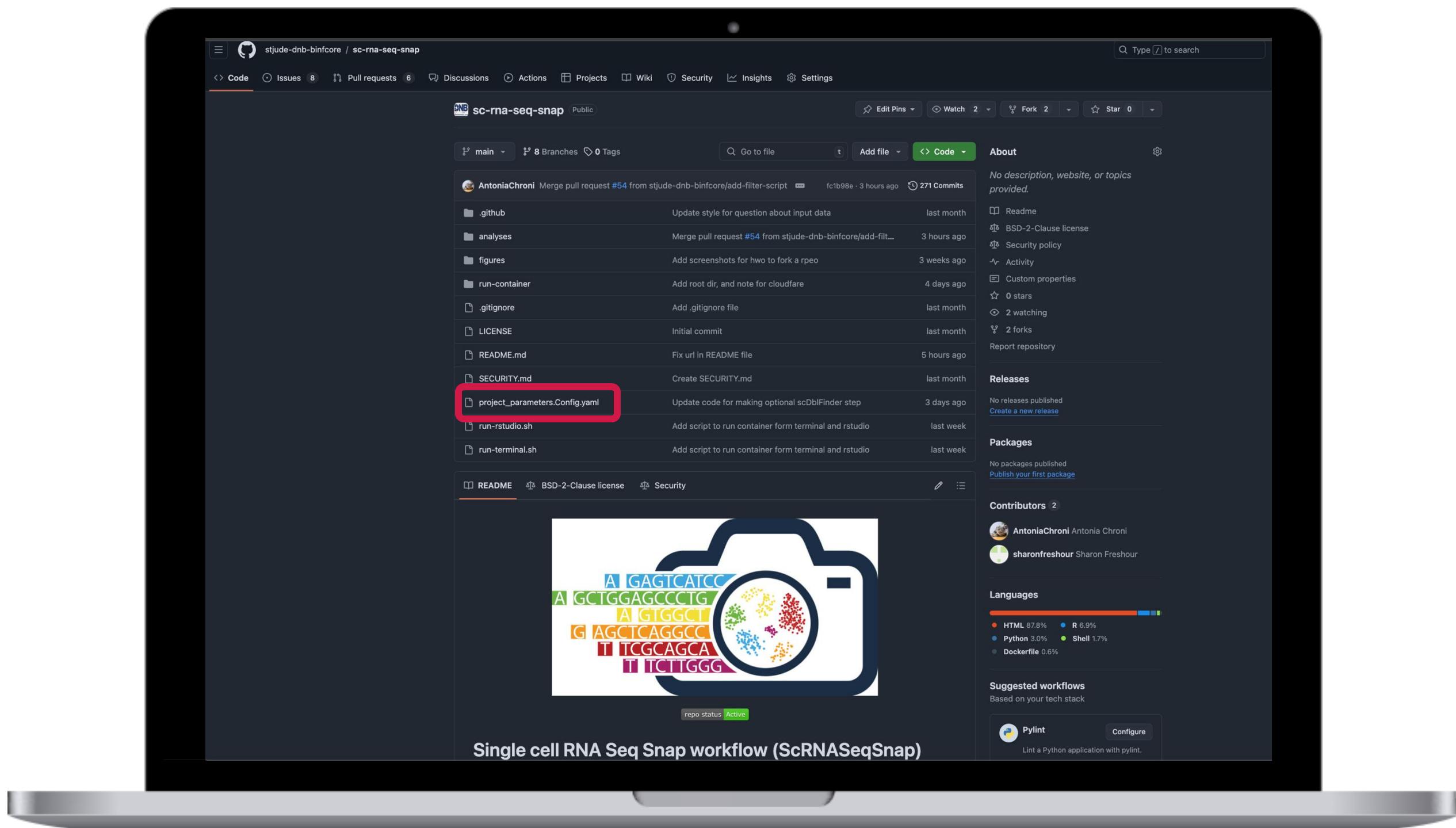
# My typical project organization system

The image shows a split-screen view illustrating a project organization system. On the left, a GitHub repository page for 'sc-rna-seq-snap' is displayed. A red arrow points from the 'analyses' folder in the repository's contents to the same folder in a file browser window on the right. The file browser window is titled 'upstream-analysis' and lists various project components:

- Downloads:
  - analyses
  - archive
  - figures
  - outline
  - project\_parameters.Config.yaml
  - analyses
  - figures
  - run-container
  - .gitignore
  - LICENSE
  - README.md
  - SECURITY.md
  - project\_parameters.Config.yaml
  - run-rstudio.sh
  - run-terminal.sh
- lectures:
  - cell-contamination-removal-analysis
  - cell-types-annotation
  - cellranger-analysis
  - cluster-cell-calling
  - data-exploratory-analysis
  - fastqc-analysis
  - integrative-analysis
  - README.md
  - upstream-analysis
  - 01A\_run\_seurat\_qc.Rmd
  - 01B\_run\_seurat\_qc\_multiple\_samples.R
  - 02\_run\_SoupX.knit.md
  - 02\_run\_SoupX.Rmd
  - 03\_run\_scDblFinder.knit.md
  - 03\_run\_scDblFinder.Rmd
  - 04\_run\_filter\_object.Rmd
  - 05\_run\_summary\_report\_v1.Rmd
  - 05\_run\_summary\_report.Rmd
  - job.err
  - job.out
  - lsf-script.txt
  - plots
  - README.md
  - Report-Final-summary-2024-11-18.log
  - Report-Final-summary-2024-11-20.log
  - Report-scDblFinder-2024-11-21.log
  - Report-seurat-qc-2024-11-01.log
  - Report-SoupX-2024-11-01.log
  - Report-SoupX-2024-11-15.log
  - results
  - run-upstream-analysis copy.R
  - run-upstream-analysis.R
  - run-upstream-analysis.sh
  - util



# Customize your project based on the experiment!



<https://github.com/stjude-dnb-bincore/sc-rna-seq-snap>





# Project parameters file

```
project_parameters.Config.yaml
```

```
1 # the following parameters are the same across the project and might be needed in more than one module #
2 root_dir: "./sc-rna-seq-snap" # path to the main dir of the project where GitHub repo lives
3 data_dir: "./sc-rna-seq-snap/analyses/cellranger-analysis/results/02_cellranger_count/ForcedCells8000Parameters" # path to data dir of the project
4 metadata_dir: "./mouse-test-dataset" # path to metadata dir of the project
5 genome_name: "GRCm39" # define genome reference and versioning
6 PROJECT_NAME: "mouse-test-dataset"
7 PI_NAME: "Stanislav Zakharenko"
8 TASK_ID: "NA"
9 PROJECT_LEAD_NAME: "NA"
10 DEPARTMENT: "Developmental Neurobiology"
11 LEAD_ANALYSTS: "Antonia Chroni, PhD"
12 GROUP_LEAD: "Cody A. Ramirez, PhD"
13 CONTACT_EMAIL: "antonia.chroni@stjude.org"
14 PIPELINE: "Standard sc-/sn-RNA-Seq Analysis in 10X Genomics data"
15 START_DATE: "10/15/2024"
16 COMPLETION_DATE: "ONGOING"
17
18
19 # the following parameters are set up as default values and/or are specific for the following modules:
20 # `./analyses/fastqc-analysis`
21 fastqc_dir: "." #path to the fastqc files for the `fastqc-analysis` module
22
23 # `./analyses/cellranger-analysis`
24 genome_reference_path: "./" #path to genome reference to be used for the `cellranger-analysis` module
25 cellranger_parameters: "ForcedCells8000Parameters" # or "ForcedCells8000Parameters" etc
26 genome_name_cellranger: "GRCm39" #please define the genome of preference for dual genomes. In case for single genomes, please use the same as above
27
28 # `./analyses/upstream-analysis`
29 print_pdf_seurat_multiple_samples: "YES" # set value for 01B_run_seurat_qc_multiple_samples.R
30 use_condition_split_seurat_multiple_samples: "NO" # set value for 01B_run_seurat_qc_multiple_samples.R
31 grouping: "orig.ident" # define grouping to use
32 Regress_Cell_Cycle_value: "NO" # Indicates whether or not to regress for cell cycle and, if so, which method to use and scale data; acceptable values: "NO" and "YES"
33 assay: "RNA" # define assay
34 min_genes: 300 # define minimum number of genes for filtering
35 min_count: 500 # define minimum number of UMIs for filtering
36 mtDNA_pct_default: 10 # define minimum percentage of mtDNA for filtering
37 normalize_method: "log_norm" # define method for normalization of counts
38 num_pcs: 30 # define number of principal components
39 nfeatures_value: 3000 # define number of variable features
40 prefix: "lognorm" # create label based on the normalization method used
41 use_miQC: "NO" # define use of miQC R package or not; see `README.md` file for more information; acceptable values: "YES" and "NO"
42 use_only_step1: "YES" # define use of both or only first step for filtering low quality cells; see `README.md` file for more information; acceptable values: "NO" and "YES"
43 condition_value: "Genotype" # define main condition of the project; this can be used for visualization purposes on the UMAPs; value to be extra
44 num_dim_seurat_qc: [20, 25] # number of PCs to use in UMAP
45 num_neighbors_seurat_qc: [10, 20, 30] # number of neighbors to use in UMAP
46 soup_fraction_value_default: 0.05 # set rho default value to use if estimated rho is > 20%
47 num_dim_filter_object: 30 # set one value for 04_run_filter_object.Rmd
48 num_neighbors_filter_object: 30 # set one value for 04_run_filter_object.Rmd
49 use_SoupX_filtering_filter_object: "YES" # set for 04_run_filter_object.Rmd
50 use_condition_split_filter_object: "YES" # set for 04_run_filter_object.Rmd
51 print_pdf_filter_object: "NO" # set for 04_run_filter_object.Rmd
52
```





# How to run the code

The screenshot shows a GitHub repository page with a dark theme. At the top, there are links for README, License, and Security, along with a file named project\_parameters.Config.yaml. A red box highlights the following text:

**To run the code in this repository:**

1. Replace the `project_parameters.Config.yaml` with your file paths and parameters.
2. Navigate to an analysis module and run the shell script of interest:

```
cd ./sc-rna-seq-snap/analyses/<module_of_interest>
```

If you have forked the repo, you will need to do the following steps before running the script of interest. If you have cloned the repo, you can skip this.

You need to do `sync fork` of your project repo at GitHub before running a module, if your branch is behind the main branch of the `stjude-dnb-binfcore/sc-rna-seq-snap:main`. This will update the main branch of your project repo with the new code and modules (if any). This will add code and not break any analyses already run in your project repo.

Then navigate to your `./sc-rna-seq-snap` project repo and ensure you are at the `main` branch. If not, you will need to `git checkout` to the main branch.

```
git branch  
git checkout main
```

Finally, `git pull` to get the most updated changes and code in your project repo.

```
git pull
```

Below is the main directory structure listing the analyses and data files used in this repository

```
└── analyses  
    ├── cellranger-analysis  
    ├── fastqc-analysis  
    └── upstream-analysis  
    ├── figures  
    ├── LICENSE  
    ├── project_parameters.Config.yaml  
    ├── README.md  
    └── SECURITY.md
```





# ScRNASeqSnap analysis modules

```
MD README.md x
Preview on Save ABC Preview Run Up Down
Source Visual
1 # How to use analysis modules in the Single cell RNA Seq Snap workflow (ScRNASeqSnap)
2
3 This repository contains a collection of analysis modules designed to process and analyze single cell and single nuclei RNA (sc/snRNA) data from 10X sequencing technology.
4
5 Each module is self-contained and can be executed independently or as part of a larger analysis pipeline. Below is a summary of each analysis module, including whether they are required or optional. Furthermore, the analysis modules should be run in the following recommended order:
6
7 1. `fastqc-analysis` module (description="Pipeline for FastQC quality control tool for high throughput sequence data analysis.", required=True)
8 2. `cellranger-analysis` module (description="Pipeline for running and summarizing Cell Ranger count for single or multiple libraries.", required=True)
9 3. `upstream-analysis` module (description="Pipeline for estimating QC metrics and filtering low quality cells.", required=True)
10 4. `integrative-analysis` module (description="Pipeline for Integrative analysis.", required=True)
11 5. `cluster-cell-calling` module (description="Pipeline for cluster cell calling and gene marker analysis.", required=True)
12 6. `cell-contamination-removal-analysis` module (description="To remove clusters and repeat steps (4) and (5), e.g. for PDX experiments.", required=False)
13 7. `cell-types-annotation` module (description="Pipeline for annotating cell types.", required=True)
14
15
16 ## Contact
17
18 Contributions, issues, and feature requests are welcome! Please feel free to check [issues](https://github.com/stjude-dnb-binfcore/sc-rna-seq-snap/issues).
19
20 ---
21
22 *These tools and pipelines have been developed by the Bioinformatic core team at the [St. Jude Children's Research Hospital](https://www.stjude.org/). These are open access materials distributed under the terms of the [BSD 2-Clause License](https://opensource.org/license/bsd-2-clause), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.*
```





# ScRNASeqSnap analysis modules

Upstream

MD README.md x

← → | Preview on Save | ABC | Preview | ⚙️ | +C Run | ↕ ↖ ↘ ↙ | ↻ ↺ ↻ ↺

Source Visual

```
1 # How to use analysis modules in the Single cell RNA Seq Snap workflow (ScRNASeqSnap)
2
3 This repository contains a collection of analysis modules designed to process and analyze single cell and single nuclei RNA (sc/snRNA) data from 10X sequencing technology.
4
5 Each module is self-contained and can be executed independently or as part of a larger analysis pipeline. Below is a summary of each analysis module, including whether they are required or optional. Furthermore, the analysis modules should be run in the following recommended order:
6
7 1. `fastqc-analysis` module (description="Pipeline for FastQC quality control tool for high throughput sequence data analysis.", required=True)
8 2. `cellranger-analysis` module (description="Pipeline for running and summarizing Cell Ranger count for single or multiple libraries.", required=True)
9 3. `upstream-analysis` module (description="Pipeline for estimating QC metrics and filtering low quality cells.", required=True)
10 4. `integrative-analysis` module (description="Pipeline for Integrative analysis.", required=True)
11 5. `cluster-cell-calling` module (description="Pipeline for cluster cell calling and gene marker analysis.", required=True)
12 6. `cell-contamination-removal-analysis` module (description="To remove clusters and repeat steps (4) and (5), e.g. for PDX experiments.", required=False)
13 7. `cell-types-annotation` module (description="Pipeline for annotating cell types.", required=True)
14
15
16 ## Contact
17
18 Contributions, issues, and feature requests are welcome! Please feel free to check [issues](https://github.com/stjude-dnb-binfcore/sc-rna-seq-snap/issues).
19
20 ---
21
22 *These tools and pipelines have been developed by the Bioinformatic core team at the [St. Jude Children's Research Hospital](https://www.stjude.org/). These are open access materials distributed under the terms of the [BSD 2-Clause License](https://opensource.org/license/bsd-2-clause), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.*
```





# `fastqc-analysis` module

description="Pipeline for FastQC quality control tool for high throughput sequence data analysis.", required=True

```
GitHub
```

results

- cell-contamination-removal-analysis
- cell-types-annotation
- cellranger-analysis
- cluster-cell-calling
- data-exploratory-analysis
- fastqc-analysis
- integrative-analysis
- upstream-analysis

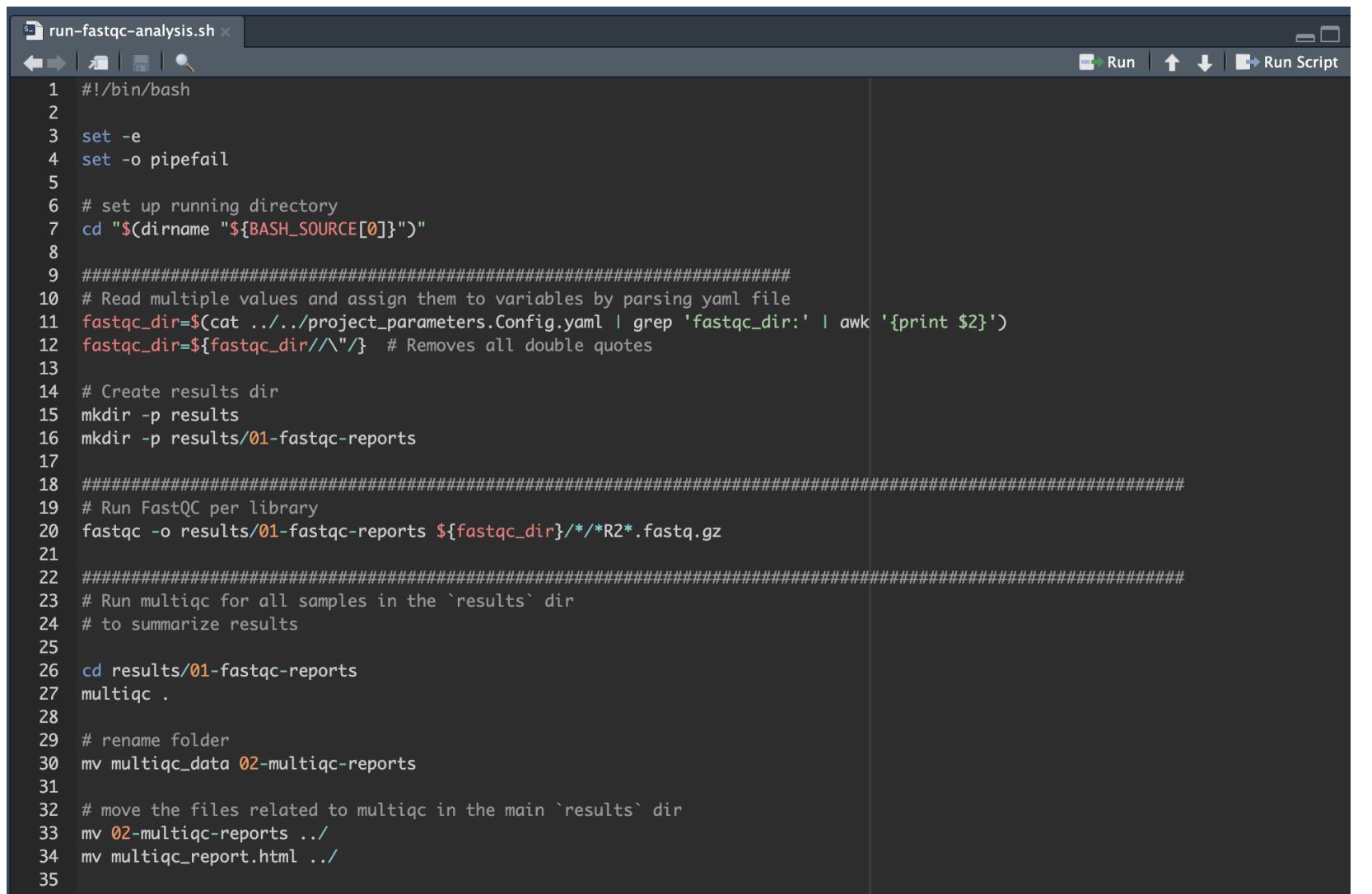
job.err  
job.out  
lsf-script.txt  
README.md  
results  
run-fastqc-analysis.sh

- 01-fastqc-reports
- 02-multiqc-reports
- fastqc
- multiqc\_report.html





# `run-fastqc-analysis.sh`



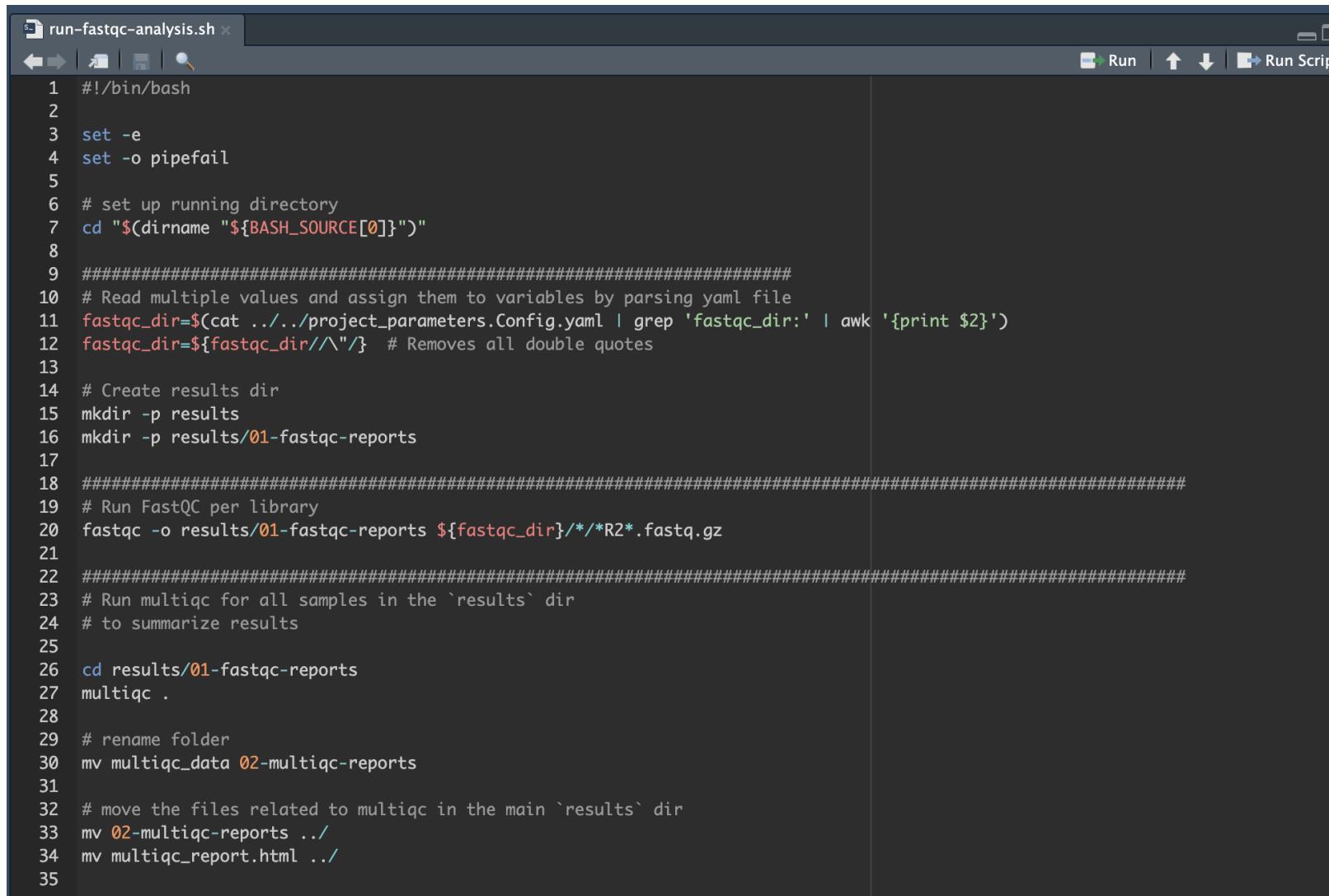
```
run-fastqc-analysis.sh x
Run | Run Script

1 #!/bin/bash
2
3 set -e
4 set -o pipefail
5
6 # set up running directory
7 cd "${dirname "${BASH_SOURCE[0]}"}"
8 #####
9 # Read multiple values and assign them to variables by parsing yaml file
10 fastqc_dir=$(cat ../../project_parameters.Config.yaml | grep 'fastqc_dir:' | awk '{print $2}')
11 fastqc_dir=${fastqc_dir//\"/} # Removes all double quotes
12
13 # Create results dir
14 mkdir -p results
15 mkdir -p results/01-fastqc-reports
16
17 #####
18 # Run FastQC per library
19 fastqc -o results/01-fastqc-reports ${fastqc_dir}/*/*R2*.fastq.gz
20
21 #####
22 # Run multiqc for all samples in the `results` dir
23 # to summarize results
24
25 cd results/01-fastqc-reports
26 multiqc .
27
28 # rename folder
29 mv multiqc_data 02-multiqc-reports
30
31 # move the files related to multiqc in the main `results` dir
32 mv 02-multiqc-reports ../
33 mv multiqc_report.html ../
34
35
```



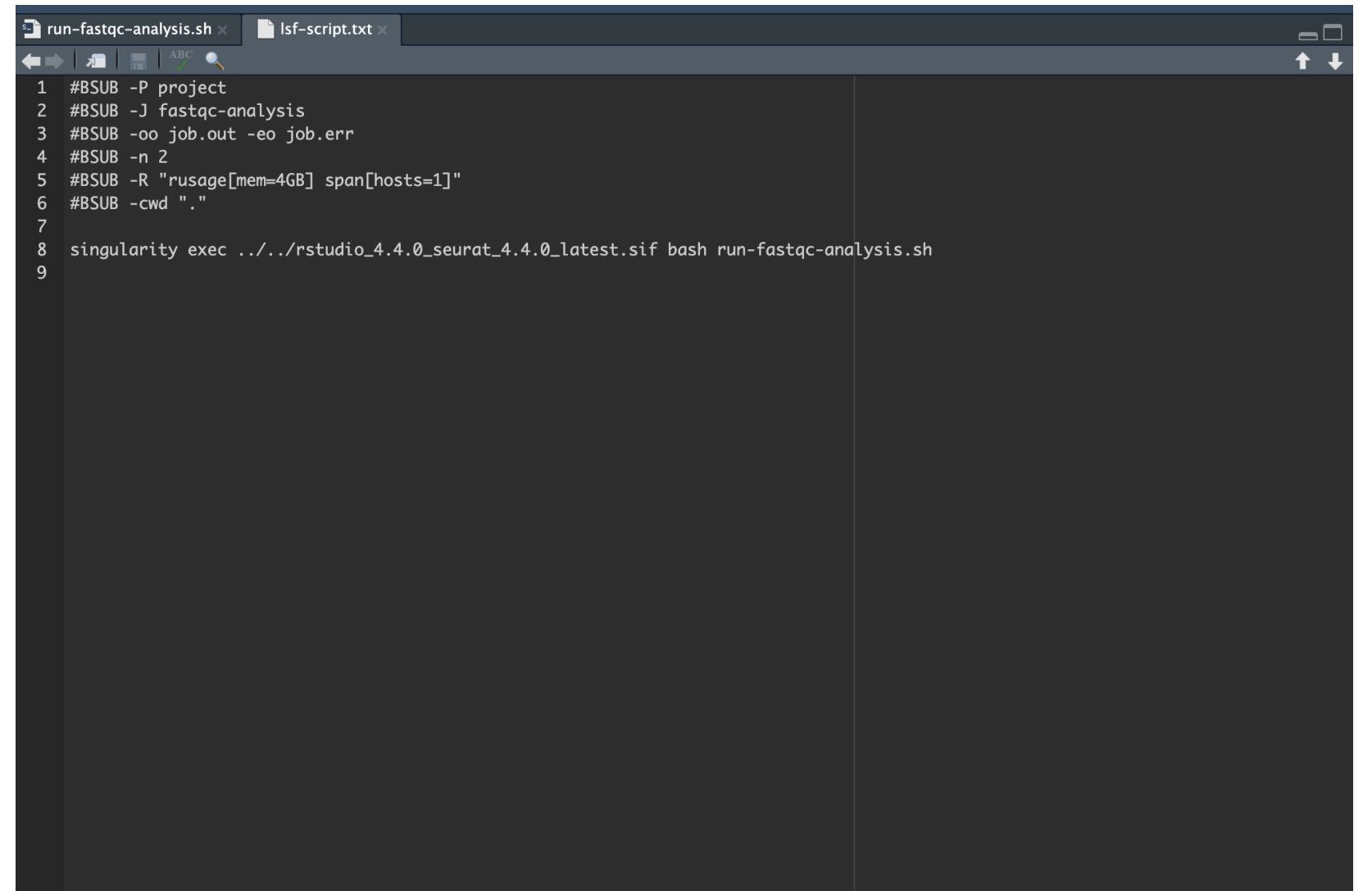


# `run-fastqc-analysis.sh` and `lsf-script.txt`



```
#!/bin/bash
set -e
set -o pipefail
# set up running directory
cd "$(dirname "${BASH_SOURCE[0]}")"
#####
# Read multiple values and assign them to variables by parsing yaml file
fastqc_dir=$(cat ../../project_parameters.Config.yaml | grep 'fastqc_dir:' | awk '{print $2}')
fastqc_dir=${fastqc_dir//"/"} # Removes all double quotes
#####
# Create results dir
mkdir -p results
mkdir -p results/01-fastqc-reports
#####
# Run FastQC per library
fastqc -o results/01-fastqc-reports ${fastqc_dir}/*/*R2*.fastq.gz
#####
# Run multiqc for all samples in the `results` dir
# to summarize results
cd results/01-fastqc-reports
multiqc .
#####
# rename folder
mv multiqc_data 02-multiqc-reports
#####
# move the files related to multiqc in the main `results` dir
mv 02-multiqc-reports ../
mv multiqc_report.html ../

```

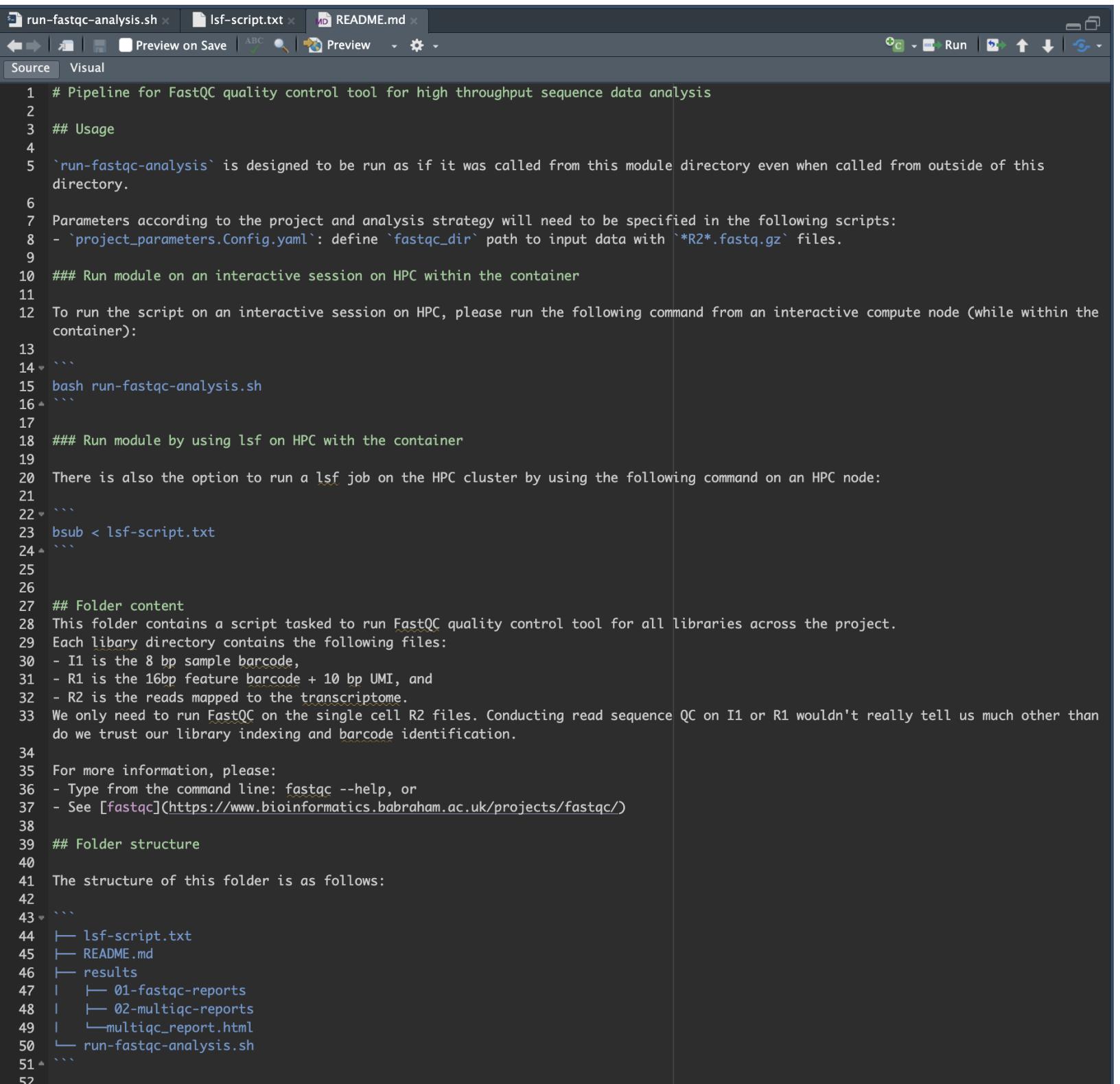


```
#BSUB -P project
#BSUB -J fastqc-analysis
#BSUB -oo job.out -eo job.err
#BSUB -n 2
#BSUB -R "rusage[mem=4GB] span[hosts=1]"
#BSUB -cwd .
singularity exec ../../rstudio_4.4.0_seurat_4.4.0_latest.sif bash run-fastqc-analysis.sh
```





# README.md



The screenshot shows a code editor window with the tab 'README.md' selected. The content of the file is as follows:

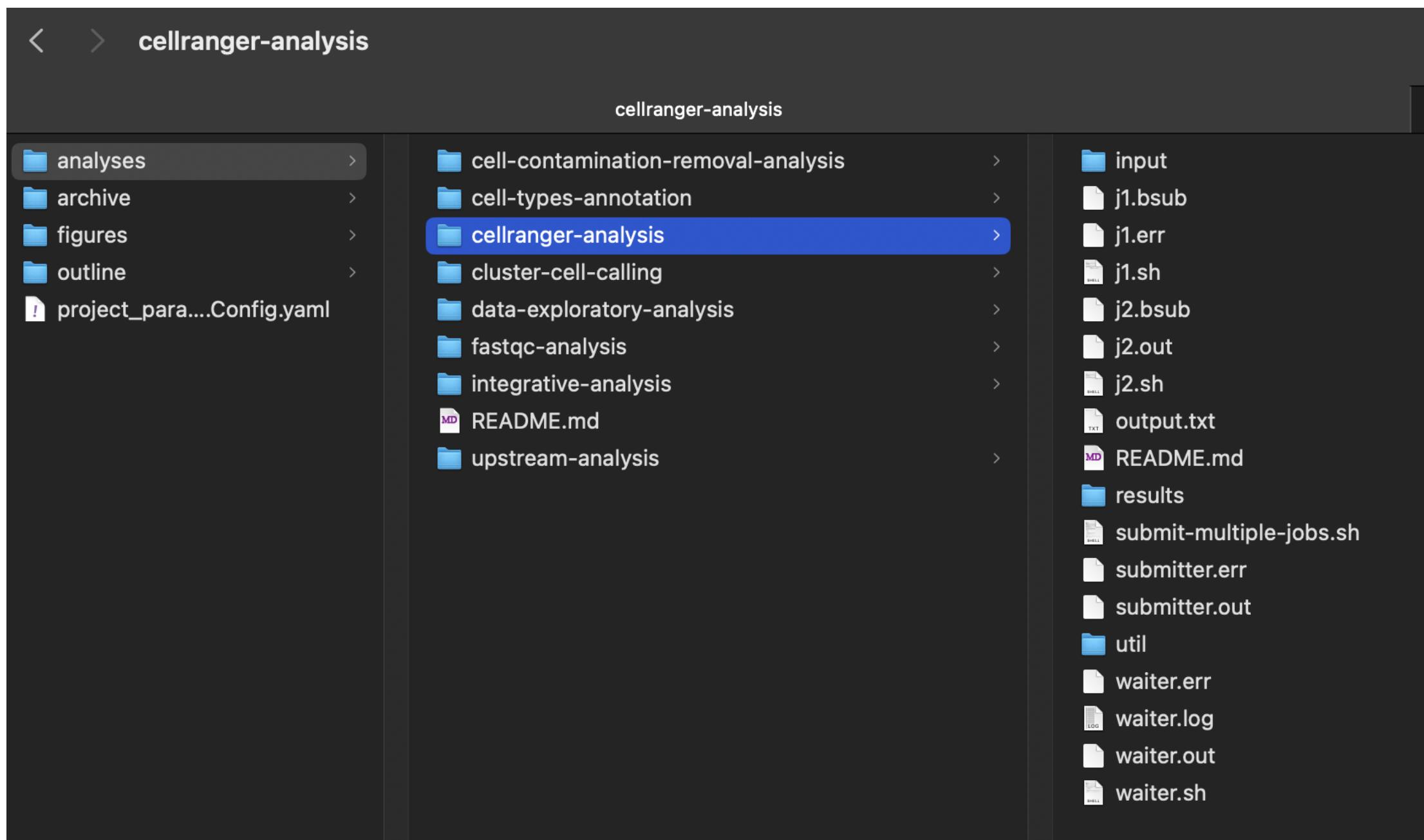
```
1 # Pipeline for FastQC quality control tool for high throughput sequence data analysis
2
3 ## Usage
4
5 `run-fastqc-analysis` is designed to be run as if it was called from this module directory even when called from outside of this
6 directory.
7
8 Parameters according to the project and analysis strategy will need to be specified in the following scripts:
9 - `project_parameters.Config.yaml` : define `fastqc_dir` path to input data with `*R2*.fastq.gz` files.
10
11 ### Run module on an interactive session on HPC within the container
12
13 To run the script on an interactive session on HPC, please run the following command from an interactive compute node (while within the
14 container):
15
16 ...
17
18 #### Run module by using lsf on HPC with the container
19
20 There is also the option to run a lsf job on the HPC cluster by using the following command on an HPC node:
21
22 ...
23 bsub < lsf-script.txt
24 ...
25
26
27 ## Folder content
28 This folder contains a script tasked to run FastQC quality control tool for all libraries across the project.
29 Each library directory contains the following files:
30 - I1 is the 8 bp sample barcode,
31 - R1 is the 16bp feature barcode + 10 bp UMI, and
32 - R2 is the reads mapped to the transcriptome.
33 We only need to run FastQC on the single cell R2 files. Conducting read sequence QC on I1 or R1 wouldn't really tell us much other than
34 do we trust our library indexing and barcode identification.
35
36 For more information, please:
37 - Type from the command line: fastqc --help, or
38 - See [fastqc](https://www.bioinformatics.babraham.ac.uk/projects/fastqc/)
39
40 ## Folder structure
41
42 The structure of this folder is as follows:
43 ...
44
45
46
47
48
49
50
51
52
```





# `cellranger-analysis` module

description="Pipeline for running and summarizing Cell Ranger count for single or multiple libraries.", required=True





# `cellranger-analysis` module

```
MD README.md x
Preview on Save ABC Preview Run Up Down
Source Visual

1 # Pipeline for running and summarizing Cell Ranger count for single or multiple libraries for sc-/sn-RNA-Seq Analysis in 10X Genomics data
2
3 ## Usage
4
5 To run all the scripts in this module from the command line sequentially, use:
6
7 ...
8 bsub < submit-multiple-jobs.sh
9 ...
10
11 The `submit-multiple-jobs.sh` script is designed to be run as if it was called from this module directory even when called from outside of this directory.
12 - Step 1: To run the `j1.sh` script to align single or multiple libraries in parallel, i.e., `run-cellranger-analysis`.
13 - Step 2: To run `j2.sh` to summarize alignment results, i.e., `summarize-cellranger-analysis`. The latter script will be on hold and executed once all libraries are aligned and `j1.sh` is complete. This is been
14 taken care of by `waiter.sh` script.
15 Parameters according to the project and analysis strategy will need to be specified in the following scripts:
16 - `../project_parameters.Config.yaml`: define the `metadata_dir`, `genome_reference_path`, `cellranger_parameters`, and `genome_name_cellranger`. Please submit an
[issue](https://github.com/stjude-dnb-binfcore/sc-rna-seq-snap/issues) to request the path to the reference genome of preference. Our team at the Bioinformatics core at DNB maintains the following genome references:
`GRCh38`, `GRCm39`, `GRCh38ANDGRCm39`, and `GRCh38_GFP_tdTomato`, for human, mouse, and dual index genomes, respectively. Otherwise, specify the path to the reference genome of your preference.
17
18 - `submit-multiple-jobs.sh`: Here, the user will need to set up the absolute path for the directory in `#BSUB -cwd` and `prefix` parameters.
19
20 - `waiter.sh`: Here, the user will need to set up the absolute path for the directory in `#BSUB -cwd` and `prefix` parameters. Also, user needs to replace `DST` with the Sample ID used for the samples of the project.
This requires that all of the Sample IDs are the same, e.g., DST800, DST802, DST811 and so on.
21
22
23 - `j1.sh`:
24   - `--file`: This is defined by default in the `./sc-rna-seq-snap/analyses/cellranger-analysis/input/` dir. If the file lives in another folder, e.g., `metadata_dir`, this should be modified accordingly in the
`../project_parameters.Config.yaml`. The file needs to be in `*.txt` file format. In addition, there is a code line to convert `*.tsv` file to `*.txt` file, if needed. If the file exists already in `./input` as
`*.txt`, this line of code will be ignored. This file can contain either one or multiple samples as long as the following minimum parameters are provided: `ID`, `SAMPLE`, and `FASTQ` columns.
25   - `--force_cells`: User can add flags as necessary for their analyses and compare alignments with CellRanger, e.g., by using `--force_cells=8000` to constraint the number of cells to be used for alignment. We
recommend to run by default, and after careful assessment to edit parameters. We have found that the default parameters set up here work well for most of the cases.
26
27 - `j2.sh`:
28   - In case of dual genomes in the experiment, user will need to add the `--genome ${genome_name_cellranger}` flag to summarize results generated by CellRanger for all libraries.
29
30
31 Default memory for running the alignment is set up to 16GB. In case more memory is needed to run a specific sample, this can be modified here:
32 - `./util/run_cellranger.py`: User can search the following and change accordingly `'-n 4 -R "rusage[mem=4000] span[hosts=1]'`, if that is necessary. However, we have found that these resources work well for most data.
33
```





# `cellranger-analysis` module

---

```
34
35 ## Folder content
36 This folder contains scripts tasked to run and summarize Cell Ranger count for single or multiple libraries for sc-/sn-RNA-Seq Analysis in 10X Genomics data across the project. For more information and updates, please
37 see [Cell Ranger support page](https://www.10xgenomics.com/support/software/cell-ranger/latest/analysis/running-pipelines/cr-gex-count).
38 This module uses CellRanger v8.0.1 for the alignment.
39
40
41 ## Folder structure
42
43 The structure of this folder is as follows:
44
45 ..
46 └── input
47   └── project_metadata.txt
48 └── j1.sh
49 └── j2.sh
50 └── README.md
51 └── results
52   ├── 01_logs
53   ├── 02_cellranger_count
54   │   └── DefaultParameters
55   └── 03_cellranger_count_summary
56 └── submit-multiple-jobs.sh
57 └── util
58   └── run_cellranger.py
59   └── summarize_cellranger_results.py
60 └── waiter.sh
61 ..
62
```





# `cellranger-analysis` module

DST800

cellranger-analysis

- input
- j1bsub
- j1.err
- j1.sh
- j2bsub
- j2.out
- j2.sh
- output.txt
- README.md
- results
  - submit-multiple-jobs.sh
  - submitter.err
  - submitter.out
  - util
  - waiter.err
  - waiter.log
  - waiter.out
  - waiter.sh

01\_logs

02\_cellranger\_count

03\_cellranger\_count\_summary

DefaultParameters

D 0

D 1

D 2

D 3

D 4

D 5

D 6

D 7

D 8

D 9

\_cmdline

\_filelist

\_finalstate

\_invocation

\_jobmode

\_log

\_mrosource

\_perf

\_perf.\_truncated\_

\_sitecheck

\_tags

\_timestamp

\_uuid

\_vdrkill

\_versions

D 0.mri.tgz

outs

SC\_RNA\_COUNTER\_CS





# `cellranger-analysis` module

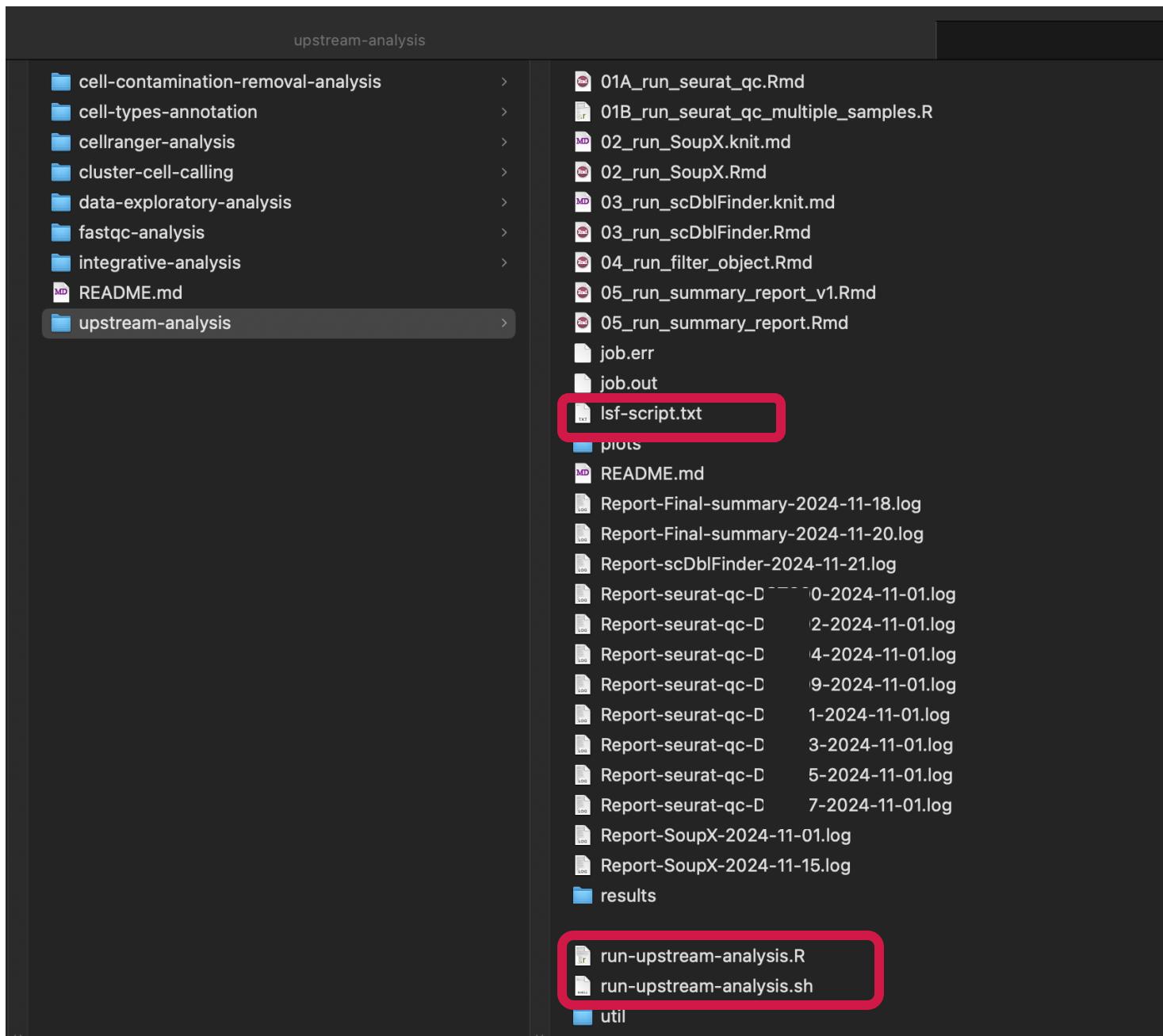
The screenshot shows a file explorer interface with two main panes. The left pane is titled '03\_cellranger\_count\_summary' and contains a list of analysis modules: cell-contamination-removal-analysis, cell-types-annotation, cellranger-analysis (which is selected), cluster-cell-calling, data-exploratory-analysis, fastqc-analysis, integrative-analysis, README.md, and upstream-analysis. The right pane is titled 'cellranger-analysis' and shows a detailed view of the 'cellranger-analysis' folder. It includes subfolders for input, j1bsub, j1.err, j1.sh, j2bsub, j2.out, j2.sh, output.txt, README.md, results, util, and waiter. A file named 'QC\_Summary\_Ce...anger\_Report.tsv' is highlighted in blue. The top bar of the file explorer has various icons for file operations like copy, paste, and search.





# `upstream-analysis` module

description="Pipeline for estimating QC metrics and filtering low quality cells.", required=True





# ‘upstream-analysis’ module

README.md

Source Visual

```
1 # Pipeline for estimating QC metrics for sc-/sn-RNA-Seq Analysis in 10X Genomics data
2
3 ## Usage
4
5 `run-upstream-analysis.sh` is designed to be run as if it was called from this module directory even when called from outside of this directory.
6
7 Parameters according to the project and analysis strategy will need to be specified in the following scripts:
8 - `project_parameters.Config.yaml` located at the `root_dir`.
9
10
11 ### Run module on an interactive session
12
13 To run all of the Rscripts in this module sequentially on an interactive session on HPC, please run the following command from an interactive compute node (while within the container):
14
15 ...
16 bash run-upstream-analysis.sh
17 ...
18
19 ### Run module by using lsf on HPC with the container
20
21 There is also the option to run a lsf job on the HPC cluster by using the following command on an HPC node:
22
23 ...
24 bsub < lsf-script.txt
25 ...
26
27
28 ## Folder content
29
30 This folder contains scripts tasked to:
31 (1) Infer QC metrics and associated plots to visually explore the quality of each library of the project.
32 (2) Evaluate QC metrics and set filters to remove low quality cells in 10x single-cell- and single-nuclei-RNA-sequencing libraries (without cell hashing experiment).
33
34 ## QC Steps and methods
35
36 The pipeline allows for the user to include/exclude methods and adjust the pipeline during QC according to the sequence type, expected cells number, type of experiment, and genome reference.
37
38 If the user does not wish to include results from step (2), this can be defined in the `project_parameters.Config.yaml` file.
39
40
41 ### (1) Seurat QC metrics
42
43 [Seurat](https://satijalab.org/seurat/articles/pbmc3k_tutorial.html) and [scooter](https://github.com/igordot/scooter) workflows are implemented to pre-process, filter and plot the RNA-sequencing data. The CellRanger output from the `cellranger-analysis` module is used for this step. User will have to define `params` as needed for their experiment.
44 - Before and after filter: Plot distribution of the number of UMI, Median genes detected per cell, Median percent mitochondrial reads per cell.
45 - Summary of Cell Statistics: Percent of reads in cells, Median UMI count per cell, Median genes detected per cell, Median percent reads mitochondrial.
46 - Data were normalized by using the global-scaling normalization method "LogNormalize" that normalizes the feature expression measurements for each cell by the total expression, multiplies this by a scale factor (10,000 by default), and log-transforms the result. Then, highly variable genes (HVGs) are selected to subset features that indicate high cell-to-cell variation in the dataset (i.e., they are highly expressed in some cells, and lowly expressed in others). Then, these HVGs are used as input to principal component analysis, and the top 30 principal components are selected. A combination of different dimensions (20, 25) and number of neighbors (30, 20, 10) are used along with the principal components to calculate the tSNE (t-Distributed Stochastic Neighbor Embedding) embeddings.
47
48 Here, the user can select to implement the following strategies to remove low quality cells:
49 - `step A` [miQC](https://bioconductor.org/packages/devel/bioc/vignettes/miQC/inst/doc/miQC.html) R package. The miQC model is based on the assumption that there are a non-trivial number of compromised cells in the dataset, which is not true in all datasets. If it is already known that the dataset is high-quality with a trivial number of compromised cells, we recommend that the user skip this step.
50 - `step B` `run_QC_default` function. This is split in two filtering steps.
51 - `step 1`: Filter cells with low content of genes expressed and remove mtDNA from each library (as defined in the `params`).
52 - `step 2`: `Find_Outlier_Threshold` function. This is an optional step (as defined in the `params`).
53 If miQC does not identify any low quality cells, then the `step B` will automatically be used as for the filtering strategy.
54
55 #### Post alignment/cell quality filtering parameters
56 We recommend that the user use the following parameters for initial QC, and then adjust accordingly if necessary:
57 - `scRNA`: min_genes = 300 (nFeature_RNA (genes detected))
58     min_count = 500 (nCount_RNA (UMIs detected))
59     mtDNA_pct_default = 15 (ideally 10; Percent Mitochondrial)
60 - `snRNA`: min_genes = 300 (nFeature_RNA (genes detected))
61     min_count = 500 (nCount_RNA (UMIs detected))
62     mtDNA_pct_default = 5 (ideally 1; Percent Mitochondrial)
```

```

63
64 #### (2) Estimating and filtering out ambient mRNA ('empty droplets')
65
66 [SoupX](https://cran.r-project.org/web/packages/SoupX/vignettes/pbmcTutorial.html) profiles "the soup", i.e., collection of cell-free mRNAs floating in the input solution. The soup looks different for each input
solution and strongly resembles the expression pattern obtained by summing all the individual cells.
67
68 SoupX calculates 'Cell-specific contamination fraction' (Estimate (or manually set) the contamination fraction, the fraction of UMIs originating from the background, in each cell) and infers a 'corrected
expression matrix' (Correct the expression of each cell using the ambient mRNA expression profile and estimated contamination).
69
70 The CellRanger output from the `cellranger-analysis` module is used for this step.
71 - Contamination summary table and Cell-specific contamination fraction plot are generated.
72
73
74 #### (3) Estimating and filtering out doublets
75
76 Popular approach of scRNAseq uses oil droplets or wells to isolate single cells along with barcoded beads. Depending on the cell density loaded, a proportion of reaction volumes (i.e. droplets or wells) will
capture more than one cell, forming 'doublets' (or 'multiplets'), i.e. two or more cells captured by a single reaction volume and thus sequenced as a single-cell artifact.
77
78 The proportion of doublets is proportional to the number of cells captured. It is common in single-cell experiments to have 10-20% doublets, making accurate doublet detection critical.
79
80 Doublets are prevalent in single-cell sequencing data and can lead to artifactual findings. We will use a computational approach to calculate and remove doublets from the library. Here, we use
[scDblFinder](https://bioconductor.org/packages/devel/bioc/vignettes/scDblFinder/inst/doc/scDblFinder.html) method for identifying doublets/multiplets in single-cell data.
81
82 The 'seurat_obj_raw.rds' object from step (1) is used for this step.
83 - Summary table with doublet metrics and doublets prediction plot are generated.
84
85
86 #### (4) Merging filtered data
87
88 Next, we merge count matrices from steps (1-3) after filtering out low quality cells, ambient RNA (optional as defined in the `params`), and doublets. Seurat object and metadata for the library along with UMAP
embeddings are saved to be used for downstream analyses.
89
90 #### (5) Final QC summary report
91
92 Lastly, we provide a final QC summary report containing graphs and summary tables across each QC step.
93
94 ## Folder structure |
95
96 The structure of this folder is as follows:
97
98 ...
99
100 |   01A_run_seurat_qc.Rmd
101 |   01B_run_seurat_qc_multiple_samples.R
102 |   02_run_SoupX.Rmd
103 |   03_run_scDblFinder.Rmd
104 |   04_run_filter_object.Rmd
105 |   05_run_summary_report.Rmd
106 |   plots
107 |   lsf-script.txt
108 |   README.md
109 |   results
110 |   run-upstream-analysis.R
111 |   run-upstream-analysis.sh
112 |   util
113 |       |   function-calculate-qc-metrics.R
114 |       |   function-create-UMAP.R
115 |       |   function-process-Seurat.R
116 |       ...

```





# `upstream-analysis` module

#	steps	aim
1	seurat-qc pipeline	To infer QC metrics and remove low quality cells (miQC R package and/or `run_QC_default`)

```
40
41 #### (1) Seurat QC metrics
42
43 [Seurat](https://satijalab.org/seurat/articles/pbmc3k\_tutorial.html) and [scooter](https://github.com/igordot/scooter) workflows are implemented to pre-process, filter and plot the RNA-sequencing data. The CellRanger output from the `cellranger-analysis` module is used for this step. User will have to define `params` as needed for their experiment.
44   - Before and after filter: Plot distribution of the number of genes, UMI, and percent mitochondrial reads per cell.
45   - Summary of Cell Statistics: Percent of reads in cells, Median UMI count per cell, Median genes detected per cell, Median percent reads mitochondrial.
46   - Data were normalized by using the global-scaling normalization method “LogNormalize” that normalizes the feature expression measurements for each cell by the total expression, multiplies this by a scale factor (10,000 by default), and log-transforms the result. Then, highly variable genes (HVGs) are selected to subset features that indicate high cell-to-cell variation in the dataset (i.e, they are highly expressed in some cells, and lowly expressed in others). Then, these HVGs are used as input to principal component analysis, and the top 30 principal components are selected. A combination of different dimensions (20, 25) and number of neighbors (30, 20, 10) are used along with the principal components to calculate the UMAP (Uniform Manifold Approximation and Projection) embeddings.
47
48 Here, the user can select to implement the following strategies to remove low quality cells:
49   - `step A` [miQC] R package. The miQC model is based on the assumption that there are a non-trivial number of compromised cells in the dataset, which is not true in all datasets. If it is already known that the dataset is high-quality with a trivial number of compromised cells, we recommend that the user skip this step.
50   - `step B` `run_QC_default` function. This is split in two filtering steps.
51     - `step 1`: Filter cells with low content of genes expressed and remove mtDNA from each library (as defined in the `params`).
52     - `step 2`: `Find_Outlier_Threshold` function. This is an optional step (as defined in the `params`).
53 If miQC does not identify any low quality cells, then the `step B` will automatically be used as for the filtering strategy.
54
55 ##### Post alignment/cell quality filtering parameters
56 We recommend that the user use the following parameters for initial QC, and then adjust accordingly if necessary:
57   - `scRNA`: min_genes = 300 (nFeature_RNA (genes detected))
58     min_count = 500 (nCount_RNA (UMIs detected))
59     mtDNA_pct_default = 15 (ideally 10; Percent Mitochondrial)
60   - `snRNA`: min_genes = 300 (nFeature_RNA (genes detected))
61     min_count = 500 (nCount_RNA (UMIs detected))
62     mtDNA_pct_default = 5 (ideally 1; Percent Mitochondrial)
63
```



# ‘upstream-analysis` module

---

#	steps	aim
1	seurat-qc pipeline	To infer QC metrics and remove low quality cells (miQC R package and/or `run_QC_default`)
2	SoupX method	To estimate contamination fraction (ambient RNA) and adjust expression matrix

```

63
64  ### (2) Estimating and filtering out ambient mRNA (`empty droplets`)
65
66 [SoupX](https://cran.r-project.org/web/packages/SoupX/vignettes/pbmcTutorial.html) profiles “the soup”, i.e., collection of cell-free mRNAs floating in the input solution. The soup looks different for each input
67 solution and strongly resembles the expression pattern obtained by summing all the individual cells.
68 SoupX calculates `Cell-specific contamination fraction` (Estimate (or manually set) the contamination fraction, the fraction of UMIs originating from the background, in each cell) and infers a `corrected
69 expression matrix` (Correct the expression of each cell using the ambient mRNA expression profile and estimated contamination).
70 The CellRanger output from the `cellranger-analysis` module is used for this step.
71 - Contamination summary table and Cell-specific contamination fraction plot are generated.
72
73

```



# ‘upstream-analysis` module

---

#	steps	aim
1	seurat-qc pipeline	To infer QC metrics and remove low quality cells (miQC R package and/or `run_QC_default`)
2	SoupX method	To estimate contamination fraction (ambient RNA) and adjust expression matrix
3	scDblFinder method	To identify doublets/multiplets

```

73
74  ### (3) Estimating and filtering out doublets
75
76 Popular approach of scRNAseq uses oil droplets or wells to isolate single cells along with barcoded beads. Depending on the cell density loaded, a proportion of reaction volumes (i.e. droplets or wells) will
77 capture more than one cell, forming ‘doublets’ (or ‘multiplets’), i.e. two or more cells captured by a single reaction volume and thus sequenced as a single-cell artifact.
78 The proportion of doublets is proportional to the number of cells captured. It is common in single-cell experiments to have 10-20% doublets, making accurate doublet detection critical.
79
80 Doublets are prevalent in single-cell sequencing data and can lead to artifactual findings. We will use a computational approach to calculate and remove doublets from the library. Here, we use
81 [ScDblFinder](https://bioconductor.org/packages/devel/bioc/vignettes/scDblFinder/inst/doc/scDblFinder.html) method for identifying doublets/multiplets in single-cell data.
82 The `seurat_obj_raw.rds` object from step (1) is used for this step.
83 - Summary table with doublet metrics and doublets prediction plot are generated.
84
85

```



# ‘upstream-analysis` module

#	steps	aim
1	seurat-qc pipeline	To infer QC metrics and remove low quality cells (miQC R package and/or `run_QC_default`)
2	SoupX method	To estimate contamination fraction (ambient RNA) and adjust expression matrix
3	scDblFinder method	To identify doublets/multiplets
4	Merge and filter data	(steps 2 and 3 are optional during final filtering)
5	Final QC summary report	To create clean object for downstream analysis

```

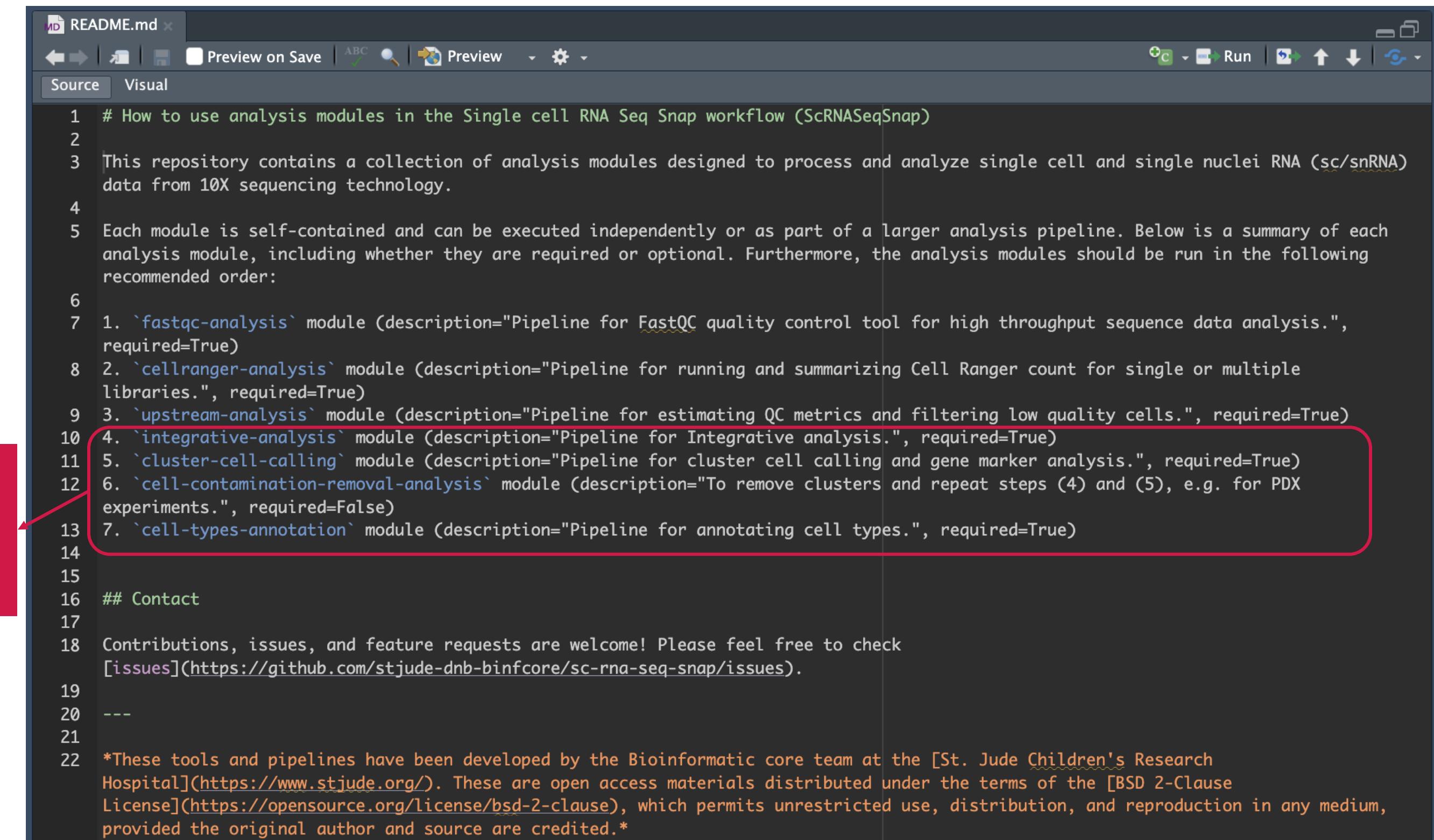
85
86  ### (4) Merging filtered data
87
88 Next, we merge count matrices from steps (1-3) after filtering out low quality cells, ambient RNA (optional as defined in the `params`), and doublets. Seurat object and metadata for the library along with UMAP
embeddings are saved to be used for downstream analyses.
89
90  ### (5) Final QC summary report
91
92 Lastly, we provide a final QC summary report containing graphs and summary tables across each QC step.
93

```





# ScRNASeqSnap analysis modules

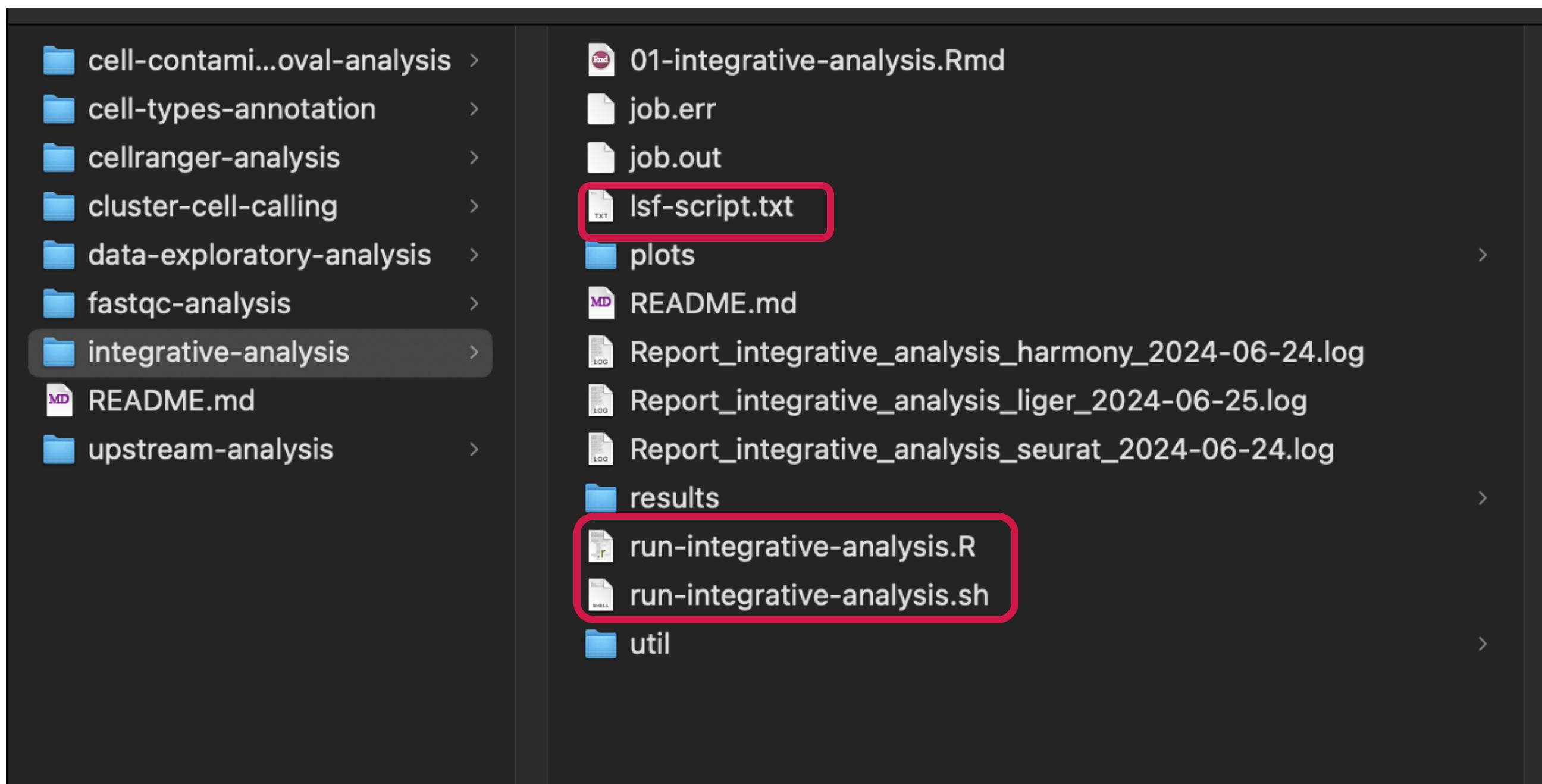


```
MD README.md x
Preview on Save ABC Preview Run Up Down
Source Visual
1 # How to use analysis modules in the Single cell RNA Seq Snap workflow (ScRNASeqSnap)
2
3 This repository contains a collection of analysis modules designed to process and analyze single cell and single nuclei RNA (sc/snRNA) data from 10X sequencing technology.
4
5 Each module is self-contained and can be executed independently or as part of a larger analysis pipeline. Below is a summary of each analysis module, including whether they are required or optional. Furthermore, the analysis modules should be run in the following recommended order:
6
7 1. `fastqc-analysis` module (description="Pipeline for FastQC quality control tool for high throughput sequence data analysis.", required=True)
8 2. `cellranger-analysis` module (description="Pipeline for running and summarizing Cell Ranger count for single or multiple libraries.", required=True)
9 3. `upstream-analysis` module (description="Pipeline for estimating QC metrics and filtering low quality cells.", required=True)
10 4. `integrative-analysis` module (description="Pipeline for Integrative analysis.", required=True)
11 5. `cluster-cell-calling` module (description="Pipeline for cluster cell calling and gene marker analysis.", required=True)
12 6. `cell-contamination-removal-analysis` module (description="To remove clusters and repeat steps (4) and (5), e.g. for PDX experiments.", required=False)
13 7. `cell-types-annotation` module (description="Pipeline for annotating cell types.", required=True)
14
15
16 ## Contact
17
18 Contributions, issues, and feature requests are welcome! Please feel free to check [issues](https://github.com/stjude-dnb-binfcore/sc-rna-seq-snap/issues).
19
20 ---
21
22 *These tools and pipelines have been developed by the Bioinformatic core team at the [St. Jude Children's Research Hospital](https://www.stjude.org/). These are open access materials distributed under the terms of the [BSD 2-Clause License](https://opensource.org/license/bsd-2-clause), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.*
```



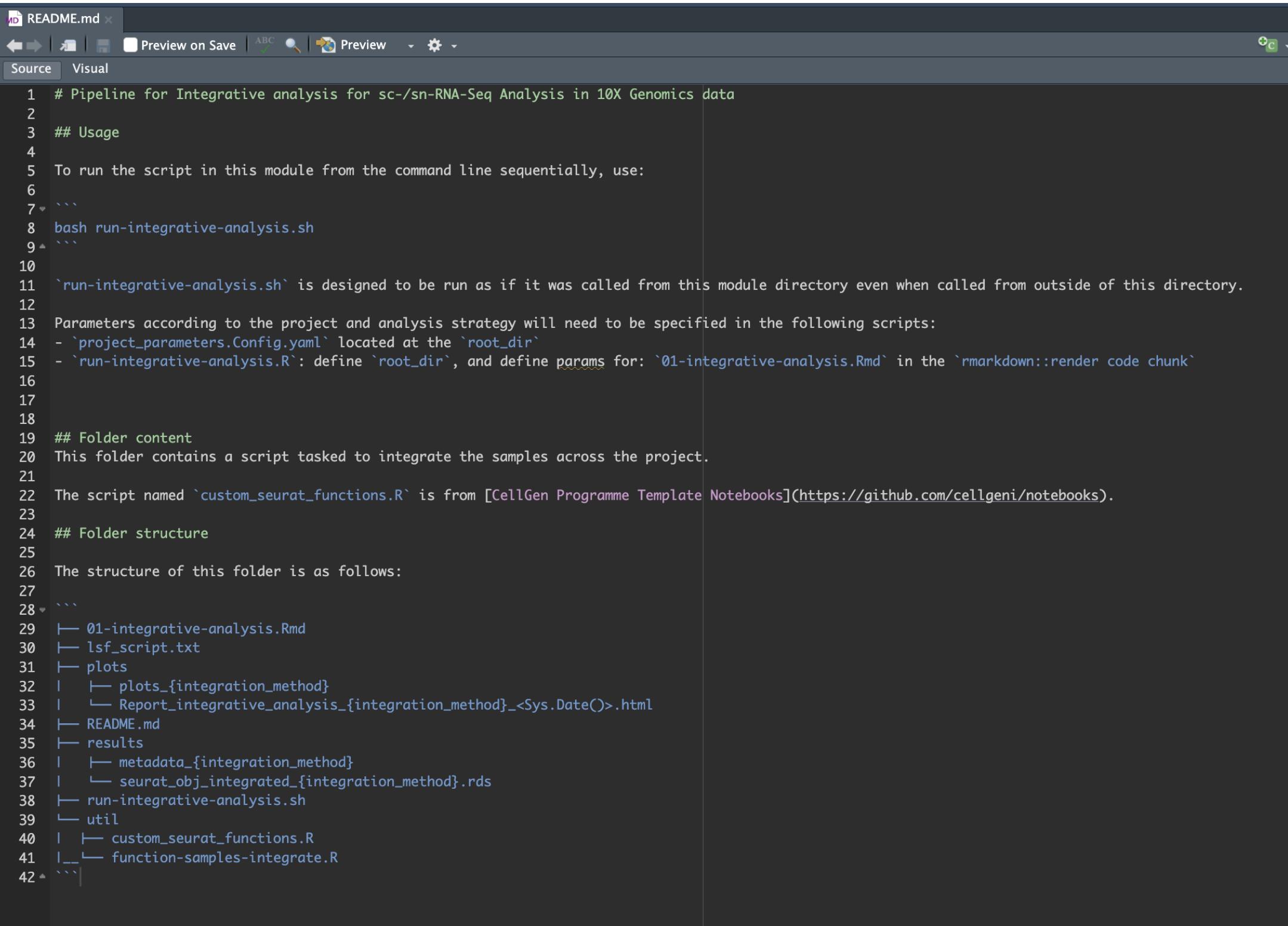
# `integrative-analysis` module

description="Pipeline for Integrative analysis.", required=True





# `integrative-analysis` module



The screenshot shows a code editor window with the file "README.md" open. The content of the file is as follows:

```
MD README.md x
Preview on Save ABC Preview Source Visual
1 # Pipeline for Integrative analysis for sc-/sn-RNA-Seq Analysis in 10X Genomics data
2
3 ## Usage
4
5 To run the script in this module from the command line sequentially, use:
6
7 ...
8 bash run-integrative-analysis.sh
9 ...
10
11 `run-integrative-analysis.sh` is designed to be run as if it was called from this module directory even when called from outside of this directory.
12
13 Parameters according to the project and analysis strategy will need to be specified in the following scripts:
14 - `project_parameters.Config.yaml` located at the `root_dir`
15 - `run-integrative-analysis.R`: define `root_dir`, and define params for: `01-integrative-analysis.Rmd` in the `rmarkdown::render` code chunk
16
17
18
19 ## Folder content
20 This folder contains a script tasked to integrate the samples across the project.
21
22 The script named `custom_seurat_functions.R` is from [CellGeni Programme Template Notebooks](https://github.com/cellgeni/notebooks).
23
24 ## Folder structure
25
26 The structure of this folder is as follows:
27
28 ...
29 |-- 01-integrative-analysis.Rmd
30 |-- lsf_script.txt
31 |-- plots
32 |   |-- plots_{integration_method}
33 |   |   Report_integrative_analysis_{integration_method}_{<Sys.Date()%>}.html
34 |-- README.md
35 |-- results
36 |   |-- metadata_{integration_method}
37 |   |   seurat_obj_integrated_{integration_method}.rds
38 |-- run-integrative-analysis.sh
39 |-- util
40 |   |-- custom_seurat_functions.R
41 |       function-samples-integrate.R
42 ...
```



# `integrative-analysis` module

#	methods	aim
1	seurat_integration	SCT method for the normalization and variance stabilization of molecular count data from scRNA-seq experiments
2	harmony_integration	Harmony algorithm projects cells into a shared embedding in which cells group by cell type rather than dataset-specific conditions. Harmony simultaneously accounts for multiple experimental and biological factors
3	liger_integration	LIGER algorithm enables the identification of shared cell types across individuals, species, and multiple modalities (gene expression, epigenetic or spatial data), as well as dataset-specific features, offering a unified analysis of heterogeneous single-cell datasets.





# `cluster-cell-calling` module

description="Pipeline for cluster cell calling and gene marker analysis.", required=True

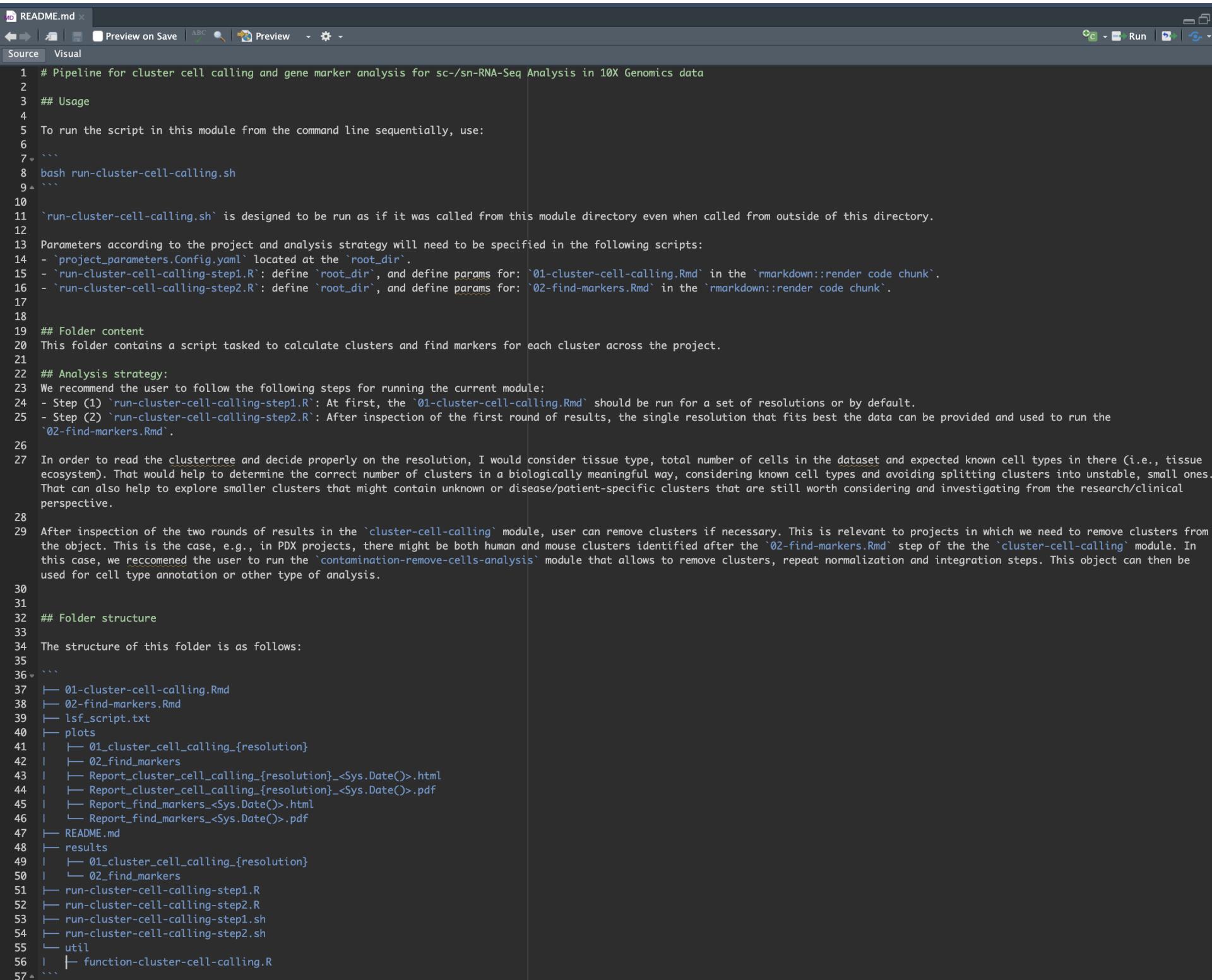
< > cluster-cell-calling

- 📁 cell-contami...oval-analysis >
- 📁 cell-types-annotation >
- 📁 cellranger-analysis >
- 📁 **cluster-cell-calling** >
  - 01-cluster-cell-calling.Rmd
  - 02-find-markers.Rmd
  - lsf-script.txt**
  - plots
  - py39
  - README.md
  - Report\_cluster\_cell\_calling\_custom\_multiple\_2024-07-02.log
  - Report\_cluster\_cell\_calling\_default\_multiple\_2024-07-02.log
  - Report\_find\_markers\_2024-07-16.log
  - results
  - run-cluster-cell-calling-step1.R
  - run-cluster-cell-calling-step1.sh
  - run-cluster-cell-calling-step2.R
  - run-cluster-cell-calling-step2.sh
  - util
- 📁 data-exploratory-analysis >
- 📁 fastqc-analysis >
- 📁 integrative-analysis >
- MD README.md
- 📁 upstream-analysis >





# `cluster-cell-calling` module



The screenshot shows a code editor window with the file "README.md" open. The content of the file is a detailed README for the "cluster-cell-calling" module. It includes instructions for running the script, details about project parameters, analysis strategy, folder content, and a folder structure diagram.

```
1 # Pipeline for cluster cell calling and gene marker analysis for sc-/sn-RNA-Seq Analysis in 10X Genomics data
2
3 ## Usage
4
5 To run the script in this module from the command line sequentially, use:
6
7 ``
8 bash run-cluster-cell-calling.sh
9 ``
10
11 `run-cluster-cell-calling.sh` is designed to be run as if it was called from this module directory even when called from outside of this directory.
12
13 Parameters according to the project and analysis strategy will need to be specified in the following scripts:
14 - `project_parameters.Config.yaml` located at the `root_dir`.
15 - `run-cluster-cell-calling-step1.R`: define `root_dir`, and define params for: `01-cluster-cell-calling.Rmd` in the `rmarkdown::render code chunk`.
16 - `run-cluster-cell-calling-step2.R`: define `root_dir`, and define params for: `02-find-markers.Rmd` in the `rmarkdown::render code chunk`.
17
18
19 ## Folder content
20 This folder contains a script tasked to calculate clusters and find markers for each cluster across the project.
21
22 ## Analysis strategy:
23 We recommend the user to follow the following steps for running the current module:
24 - Step (1) `run-cluster-cell-calling-step1.R`: At first, the `01-cluster-cell-calling.Rmd` should be run for a set of resolutions or by default.
25 - Step (2) `run-cluster-cell-calling-step2.R`: After inspection of the first round of results, the single resolution that fits best the data can be provided and used to run the `02-find-markers.Rmd`.
26
27 In order to read the cluster tree and decide properly on the resolution, I would consider tissue type, total number of cells in the dataset and expected known cell types in there (i.e., tissue ecosystem). That would help to determine the correct number of clusters in a biologically meaningful way, considering known cell types and avoiding splitting clusters into unstable, small ones. That can also help to explore smaller clusters that might contain unknown or disease/patient-specific clusters that are still worth considering and investigating from the research/clinical perspective.
28
29 After inspection of the two rounds of results in the `cluster-cell-calling` module, user can remove clusters if necessary. This is relevant to projects in which we need to remove clusters from the object. This is the case, e.g., in PDX projects, there might be both human and mouse clusters identified after the `02-find-markers.Rmd` step of the the `cluster-cell-calling` module. In this case, we reccomend the user to run the `contamination-remove-cells-analysis` module that allows to remove clusters, repeat normalization and integration steps. This object can then be used for cell type annotation or other type of analysis.
30
31
32 ## Folder structure
33
34 The structure of this folder is as follows:
35
36 ``
37 └── 01-cluster-cell-calling.Rmd
38 └── 02-find-markers.Rmd
39 └── lsf_script.txt
40 └── plots
41   └── 01_cluster_cell_calling_{resolution}
42   └── 02_find_markers
43   └── Report_cluster_cell_calling_{resolution}_{<Sys.Date()>.html}
44   └── Report_cluster_cell_calling_{resolution}_{<Sys.Date()>.pdf}
45   └── Report_find_markers_{<Sys.Date()>.html}
46   └── Report_find_markers_{<Sys.Date()>.pdf}
47 └── README.md
48 └── results
49   └── 01_cluster_cell_calling_{resolution}
50   └── 02_find_markers
51 └── run-cluster-cell-calling-step1.R
52 └── run-cluster-cell-calling-step2.R
53 └── run-cluster-cell-calling-step1.sh
54 └── run-cluster-cell-calling-step2.sh
55 └── util
56   └── function-cluster-cell-calling.R
57 ``
```





# `cluster-cell-calling` module

---

#	steps	aim
1	cluster-cell-calling	To cluster the cells under different resolutions to define granularity
2	find-markers	To find differentially expressed features (cluster biomarkers)





# `cell-contamination-removal-analysis` module - optional

description="To remove clusters and repeat steps (4) and (5), e.g. for PDX experiments.", required=False

cell-contamination-removal-analysis

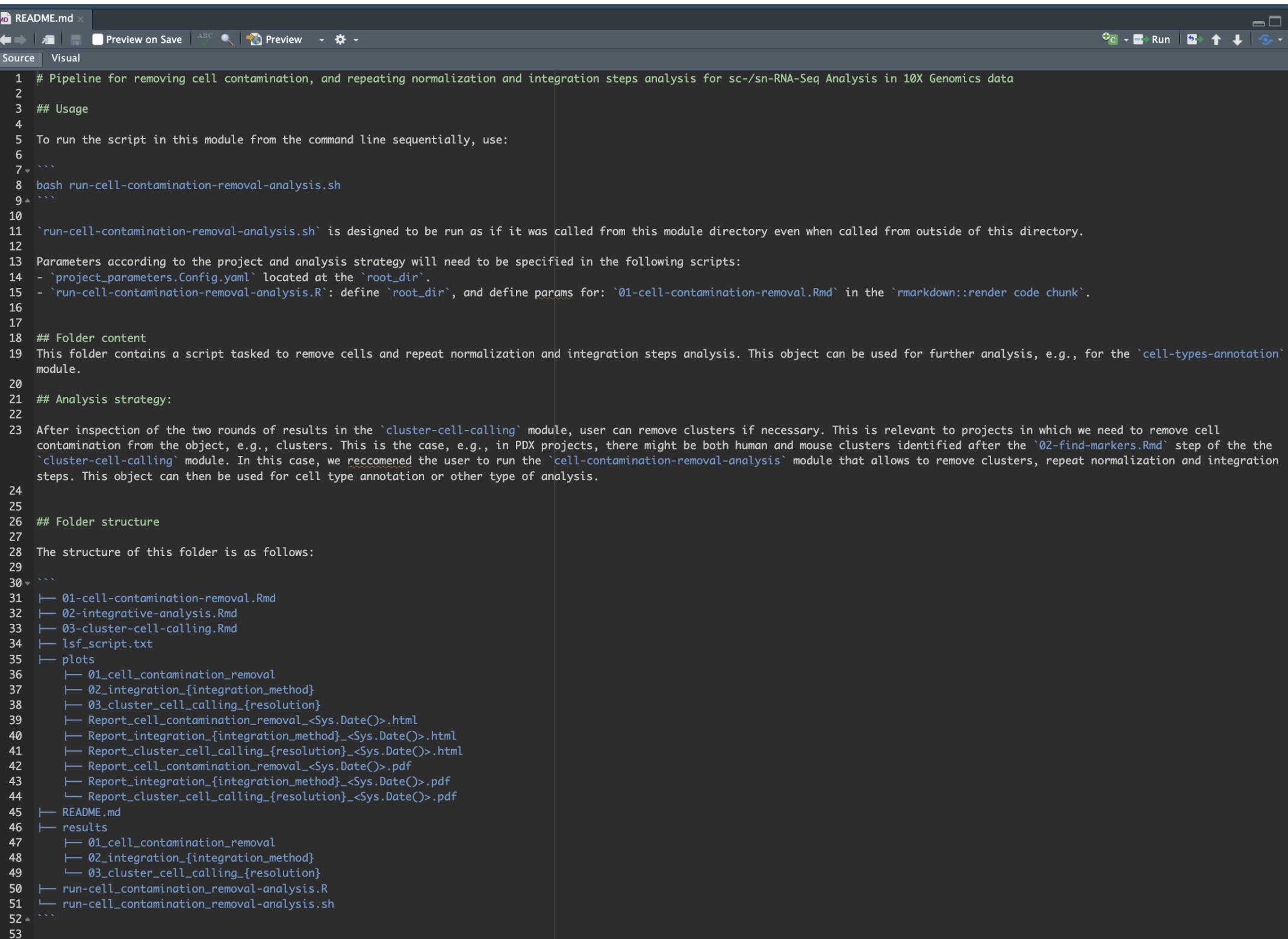
- cell-contamination-removal-analysis
- cell-types-annotation
- cellranger-analysis
- cluster-cell-calling
- data-exploratory-analysis
- fastqc-analysis
- integrative-analysis
- README.md
- upstream-analysis

- 01-cell-contamination-removal.Rmd
- 02-integrative-analysis.Rmd
- 03-cluster-cell-calling.Rmd
- 04-find-markers.Rmd
- job.err
- job.out
- lsf-script.txt
- plots
- py39
- README.md
- Report\_cell\_contamination\_removal\_2024-07-17.log
- Report\_cluster\_cell\_calling\_custom\_multiple\_2024-07-17.log
- Report\_find\_markers\_2024-07-23.log
- Report\_integrative\_analysis\_harmony\_2024-07-17.log
- results
- run-cell-contamination-removal-analysis-all-steps.R
- run-cell-contamination-removal-analysis.R
- run-cell-contamination-removal-analysis.sh





# `cell-contamination-removal-analysis` module - optional



The screenshot shows a code editor window with the title "README.md". The content of the file is a Markdown document providing instructions for using the module. It includes sections on usage, parameters, folder content, analysis strategy, and folder structure. The text is as follows:

```
1 # Pipeline for removing cell contamination, and repeating normalization and integration steps analysis for sc-/sn-RNA-Seq Analysis in 10X Genomics data
2
3 ## Usage
4
5 To run the script in this module from the command line sequentially, use:
6
7 ``
8 bash run-cell-contamination-removal-analysis.sh
9 ``
10
11 `run-cell-contamination-removal-analysis.sh` is designed to be run as if it was called from this module directory even when called from outside of this directory.
12
13 Parameters according to the project and analysis strategy will need to be specified in the following scripts:
14 - `project_parameters.Config.yaml` located at the `root_dir`.
15 - `run-cell-contamination-removal-analysis.R`: define `root_dir`, and define params for: `01-cell-contamination-removal.Rmd` in the `rmarkdown::render` code chunk.
16
17
18 ## Folder content
19 This folder contains a script tasked to remove cells and repeat normalization and integration steps analysis. This object can be used for further analysis, e.g., for the `cell-types-annotation` module.
20
21 ## Analysis strategy:
22
23 After inspection of the two rounds of results in the `cluster-cell-calling` module, user can remove clusters if necessary. This is relevant to projects in which we need to remove cell contamination from the object, e.g., clusters. This is the case, e.g., in PDX projects, there might be both human and mouse clusters identified after the `02-find-markers.Rmd` step of the the `cluster-cell-calling` module. In this case, we reccomened the user to run the `cell-contamination-removal-analysis` module that allows to remove clusters, repeat normalization and integration steps. This object can then be used for cell type annotation or other type of analysis.
24
25
26 ## Folder structure
27
28 The structure of this folder is as follows:
29
30 ``
31 ├── 01-cell-contamination-removal.Rmd
32 ├── 02-integrative-analysis.Rmd
33 ├── 03-cluster-cell-calling.Rmd
34 └── lsf_script.txt
35 └── plots
36     ├── 01_cell_contamination_removal
37     ├── 02_integration_{integration_method}
38     ├── 03_cluster_cell_calling_{resolution}
39     ├── Report_cell_contamination_removal_{Sys.Date()}.html
40     ├── Report_integration_{integration_method}_{Sys.Date()}.html
41     ├── Report_cluster_cell_calling_{resolution}_{Sys.Date()}.html
42     ├── Report_cell_contamination_removal_{Sys.Date()}.pdf
43     ├── Report_integration_{integration_method}_{Sys.Date()}.pdf
44     └── Report_cluster_cell_calling_{resolution}_{Sys.Date()}.pdf
45 └── README.md
46 └── results
47     ├── 01_cell_contamination_removal
48     ├── 02_integration_{integration_method}
49     └── 03_cluster_cell_calling_{resolution}
50 └── run-cell-contamination-removal-analysis.R
51 └── run-cell-contamination-removal-analysis.sh
52 ``
53 ``
```





# `cell-types-annotation` module

description="Pipeline for annotating cell types.", required=True

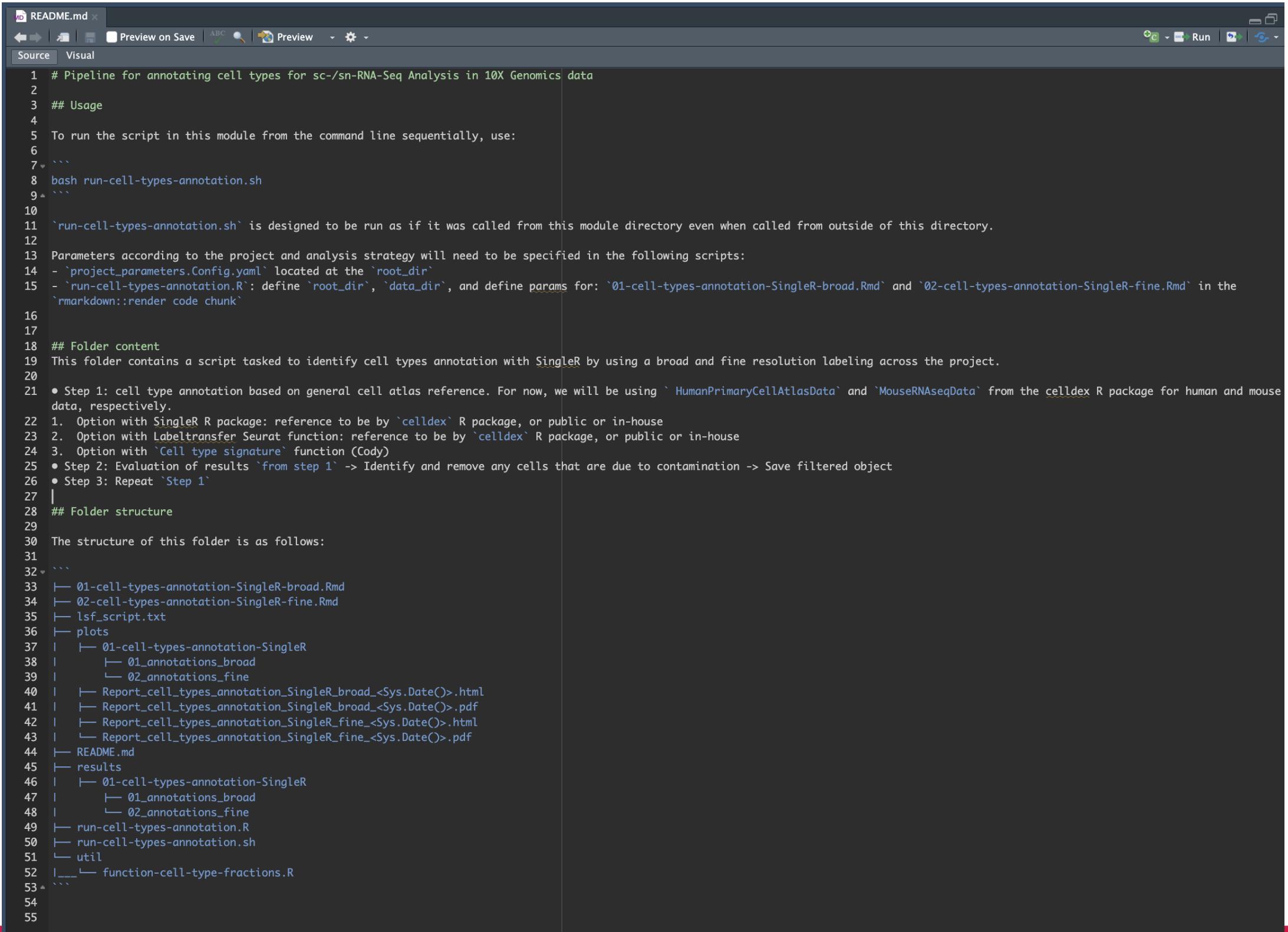
< > cell-types-annotation

- 📁 cell-contamination-removal-analysis >
- 📁 cell-types-annotation > **lsf-script.txt**
- 📁 cellranger-analysis >
- 📁 cluster-cell-calling >
- 📁 data-exploratory-analysis >
- 📁 fastqc-analysis >
- 📁 integrative-analysis >
- MD README.md
- 📁 upstream-analysis >
- 01-cell-types-annotation-SingleR-broad.Rmd
- 02-cell-types-annotation-SingleR-fine.Rmd
- plots
- README.md
- Report\_cell\_types\_annotation\_SingleR\_broad\_2024-07-17.log
- Report\_cell\_types\_annotation\_SingleR\_fine\_2024-07-17.log
- results
  - run-cell-types-annotation.R
  - run-cell-types-annotation.sh
- util





# `cell-types-annotation` module



The screenshot shows a code editor window with the file "README.md" open. The content of the file is as follows:

```
1 # Pipeline for annotating cell types for sc-/sn-RNA-Seq Analysis in 10X Genomics data
2
3 ## Usage
4
5 To run the script in this module from the command line sequentially, use:
6 ...
7 ...
8 bash run-cell-types-annotation.sh
9 ...
10
11 `run-cell-types-annotation.sh` is designed to be run as if it was called from this module directory even when called from outside of this directory.
12
13 Parameters according to the project and analysis strategy will need to be specified in the following scripts:
14 - `project_parameters.Config.yaml` located at the `root_dir`
15 - `run-cell-types-annotation.R`: define `root_dir`, `data_dir`, and define params for: `01-cell-types-annotation-SingleR-broad.Rmd` and `02-cell-types-annotation-SingleR-fine.Rmd` in the
`rmarkdown::render` code chunk
16
17
18 ## Folder content
19 This folder contains a script tasked to identify cell types annotation with SingleR by using a broad and fine resolution labeling across the project.
20
21 • Step 1: cell type annotation based on general cell atlas reference. For now, we will be using `HumanPrimaryCellAtlasData` and `MouseRNAseqData` from the celldex R package for human and mouse
data, respectively.
22 1. Option with SingleR R package: reference to be by `celldex` R package, or public or in-house
23 2. Option with Labeltransfer Seurat function: reference to be by `celldex` R package, or public or in-house
24 3. Option with `Cell type signature` function (Cody)
25 • Step 2: Evaluation of results `from step 1` -> Identify and remove any cells that are due to contamination -> Save filtered object
26 • Step 3: Repeat `Step 1`
27 |
28 ## Folder structure
29
30 The structure of this folder is as follows:
31 ...
32 ...
33 |__ 01-cell-types-annotation-SingleR-broad.Rmd
34 |__ 02-cell-types-annotation-SingleR-fine.Rmd
35 |__ lsf_script.txt
36 |__ plots
37 |   |__ 01-cell-types-annotation-SingleR
38 |   |   |__ 01_annotations_broad
39 |   |   |__ 02_annotations_fine
40 |   |   |__ Report_cell_types_annotation_SingleR_broad_<Sys.Date()>.html
41 |   |   |__ Report_cell_types_annotation_SingleR_broad_<Sys.Date()>.pdf
42 |   |   |__ Report_cell_types_annotation_SingleR_fine_<Sys.Date()>.html
43 |   |   |__ Report_cell_types_annotation_SingleR_fine_<Sys.Date()>.pdf
44 |__ README.md
45 |__ results
46 |   |__ 01-cell-types-annotation-SingleR
47 |   |   |__ 01_annotations_broad
48 |   |   |__ 02_annotations_fine
49 |__ run-cell-types-annotation.R
50 |__ run-cell-types-annotation.sh
51 |__ util
52 |___ function-cell-type-fractions.R
53 ...
54 ...
55 ...
```



## ‘cell-types-annotation` module

---

#	steps	aim
1	cell-types-annotation-SingleR-broad	To identify cell types by using SingleR and celldex database and run diagnostics
2	cell-types-annotation-SingleR-fine	To identify cell types by using SingleR and celldex database and run diagnostics





## Next steps

---

- Pipeline development (and testing phase)
- Testing phase and code review process**
- Validation phase by using other datasets/experiments
- Launching!





# Future analysis modules

---

- `clone-phylogeny-analysis` module
- `infer-tumor-microenvironment-analysis` module
- `differential-expression-analysis` module
- `trajectories-pseudotime-analysis` module
- `rna-velocity-analysis` module
- `cell-types-annotation` module: using custom atlas reference or gene signature score
- R shiny app



## More resources



- [Automation and Reproducibility in Computational Biology](#)
- [Documentation and Usability](#)
- [GitHub Docs](#)
- [Docker for beginners](#)
- [Docker for Data Scientists](#)

[./resources/links-to-resources.md](#)



# Ask me questions!



Tonia  
**Antonia Chroni, PhD**  
she/her/hers

Senior Bioinformatics Research Scientist  
New York, New York

 @AntoniaChroni



DEMO

